

UNIVERSIDADE FEDERAL DO PARANÁ

JOSUÉ OLIVEIRA CAMARGO

ANÁLISE FILOGENÉTICA DA REGIÃO CODIFICANTE DO GENOMA PLASTIDIAL:  
UMA ABORDAGEM LIVRE DE ALINHAMENTO UTILIZANDO O ALGORITMO SVECT

CURITIBA

2018

JOSUÉ OLIVEIRA CAMARGO

ANÁLISE FILOGENÉTICA DA REGIÃO CODIFICANTE DO GENOMA PLASTIDIAL:  
UMA ABORDAGEM LIVRE DE ALINHAMENTO UTILIZANDO O ALGORITMO SVECT

Dissertação apresentada como requisito parcial para obtenção do grau de Mestre pelo Programa de Pós-Graduação em Bioinformática, Setor de Educação Profissional e Tecnológica, Universidade Federal do Paraná, área de concentração Bioinformática.

Orientador: Profº Drº Fabio de Oliveira Pedrosa

Co-orientador: Profº Drº Roberto Tadeu Raittz

CURITIBA

2018

Catálogo na publicação  
Sistema de Bibliotecas UFPR  
Biblioteca de Educação Profissional e Tecnológica

Camargo, Josué Oliveira  
C172      Análise filogenética da região codificante do genoma plastidial: uma  
abordagem livre de alinhamento utilizando o algoritmo SVect / Josué Oliveira  
Camargo. - Curitiba, 2018.  
103 p.: il., tabs, grafs.

Orientador: Fabio de Oliveira Pedrosa  
Co-orientador: Roberto Tadeu Raittz  
Dissertação (Mestrado) – Universidade Federal do Paraná, Setor de  
Educação Profissional e Tecnológica, Curso de Pós-Graduação em  
Bioinformática.

1. Filogenia. 2. Plastídio. 3. Bioinformática. I. Pedrosa, Fabio de Oliveira.  
II. Raittz, Roberto Tadeu. III. Título. IV. Universidade Federal do Paraná.

CDD 576.88



MINISTÉRIO DA EDUCAÇÃO  
UNIVERSIDADE FEDERAL DO PARANÁ  
SETOR DE EDUCAÇÃO PROFISSIONAL E TECNOLÓGICA

Pós-Graduação em Bioinformática WWW.BIOINFO.UFPR.BR  
E-mail: bioinfo@ufpr.br Tel: 41 33614906

### TERMO DE APROVAÇÃO

Os membros da Banca Examinadora designada pelo Colegiado do Programa de Pós-Graduação em BIOINFORMÁTICA da Universidade Federal do Paraná foram convocados para realizar a arguição da Dissertação de Mestrado de **JOSUÉ OLIVEIRA CAMARGO** intitulada: **"ANÁLISE FILOGENÉTICA DA REGIÃO CODIFICANTE DO GENOMA PLASTIDIAL: UMA ABORDAGEM LIVRE DE ALINHAMENTO UTILIZANDO O ALGORITMO SVECT"**, após terem inquirido o aluno e realizado a avaliação do trabalho, são de parecer pela sua aprovação no rito de defesa.

A outorga do título de mestre está sujeita à homologação pelo colegiado, ao atendimento de todas as indicações e correções solicitadas pela banca e ao pleno atendimento das demandas regimentais do Programa de Pós-Graduação.

Curitiba, 27 de março de 2018.

Dr. Fábio de Oliveira Pedrosa  
Presidente  
Programa de Pós-graduação em Bioinformática – UFPR  
Departamento de Bioquímica UFPR

Dr.ª Leila do Nascimento Vieira  
Avaliadora Externa  
Pós-doutoranda da Universidade Federal de Santa Catarina - UFSC

Dr. Diéval Guizelini  
Avaliador Interno  
Programa de Pós-graduação em Bioinformática - UFPR



À minha mãe Eloisa,  
À minha esposa Joyce,  
Ao meu irmão Jhonatan,  
E a todos aqueles que acreditaram em mim.

## AGRADECIMENTOS

Agradeço primeiramente ao eterno e bom Deus por mais essa etapa concluída em minha jornada... por ter me concedido vida, força de vontade, paciência, resiliência e persistência. Obrigado por guiar meus passos.

Agradeço à toda a minha família, e em especial minha mãe, Eloísa Helena de Oliveira Camargo, por tudo o que fez por mim. Por não ter deixado o meu sonho morrer, apesar das dificuldades e obstáculos pelo caminho. Obrigado por ter me fornecido meios para aqui chegar. E se aqui cheguei, tenha a certeza que foi principalmente por você.

Agradeço ainda ao meu irmão, Jhonatan Oliveira Camargo, por todo apoio em meus piores momentos. Creio que você me conhece de verdade, pois estava lá quando eu mais precisei. E sim, se aqui cheguei, foi porque você não me deixou desistir também.

Agradeço à minha esposa, Joyce Thaise Ramos de Souza Camargo, simplesmente por tudo. Da compreensão dispendida por todo tempo que deixei de passar contigo para dispende na execução e na escrita deste trabalho, até os puxões de orelha quando começava a me deixar levar pela preguiça. Se hoje me tornei um mestre, o mérito é seu também. Te amo muito!

Agradeço ao Clube de Desbravadores e a Igreja Adventista do Sétimo dia, por ter lançado os meus fundamentos morais, éticos, de respeito à natureza e de reverência a Deus. Se hoje estou em cima, é devido também a vocês.

Agradeço aos meus amigos Marcos, André, Carlos, Juliar e Italo, por me acompanharem na jornada da vida.

Agradeço aos meus orientadores, os professores Fábio de Oliveira Pedrosa e Roberto Tadeu Raittz, por me mostrarem o caminho por trilhar a vida acadêmica. Grandes amigos que levarei para sempre comigo.

Agradeço aos mestres, companheiros de turma Willian, Daniel, Bruno, Alan e Matheus, e em especial, Dionata, Kelin, e Aryel por todo apoio, discussões, experiências vivenciadas e descontrações a caminho do RU. Vocês marcaram minha vida.

Agradeço a todos os integrantes do grupo de pesquisa de inteligência artificial aplicada à bioinformática, Bruno, Camila, Letícia, Mariane, Amanda e Diogo por me acolherem e me auxiliarem da forma que fizeram.

Agradeço à secretária do Programa de Pós-Graduação em Bioinformática da UFPR, Suzana de Azevedo Gobetti, por sua amizade, companheirismo, profissionalismo e eficiência.

Agradeço ao Programa de Pós-Graduação em Bioinformática da UFPR. Órgãos de fomento, colegiado do curso, colegas, funcionários, professores, enfim, a todos que de algum modo me auxiliaram, mesmo anonimamente, meu muito obrigado.

“É quando nos esquecemos de nós mesmos que  
fazemos coisas que jamais serão esquecidas”

Ellen G. White

## RESUMO

Com o advento do sequenciamento massivo paralelo de DNA, o campo da bioinformática tornou-se desafiador no que diz respeito a análises genômicas envolvendo *Big Data*. A filogenômica possui duas principais abordagens que utilizam os dados em larga escala para gerar reconstruções filogenéticas: super-árvores e super-matrizes. Para análises filogenéticas em geral, a limitação dos métodos baseados em alinhamento múltiplo para a comparação de sequências se torna evidente devido à complexidade algorítmica para se obter um alinhamento. Para tornar o método viável, poucas sequências são utilizadas para um grande número de organismos ou vice-versa. Podemos perder informação evolutiva ao usarmos um número mínimo de sequências, ou o prejuízo pode ser ocasionado no tempo computacional. Há a necessidade de um método capaz de gerar filogenias consistentes a partir de grandes volumes de dados em tempo computacional não-proibitivo. A partir do SVect, criamos uma metodologia livre de alinhamento para análises filogenéticas em larga escala utilizando a redução de dimensionalidade. O método se mostrou rápido, eficiente, conciso e capaz de analisar um grande volume de dados, inviável para o alinhamento múltiplo de sequências. Desenvolvemos a primeira árvore filogenética global de organismos portadores de genoma plastidial utilizando as regiões codificantes dos plastomas, e pretendemos disponibilizar tanto a árvore como a ferramenta de análise para uso científico.

Palavras-chave: Filogenia. Região codificante. SVect. Plastídio.

## **ABSTRACT**

*With the advent of large-scale genomic sequencing, the field of bioinformatics has become challenging with respect to genomic and proteomic analyzes involving Big Data. The phylogenomics have two main approaches that use the data in large scale to generate phylogenetic reconstructions: super-trees and super-matrices. For phylogenetic analyzes in general, the limitation of multiple alignment based methods for sequence comparison becomes apparent because of the algorithmic complexity to obtain an alignment. To make the method feasible, few sequences are used for a large number of organisms or vice versa. We can lose evolutionary information by using a minimal number of sequences, or the damage can be caused by computational time. Has the need for a method capable of generating consistent phylogenies from large data volumes in non-prohibitive computational time. From SVect, we created a free alignment methodology for large-scale phylogenetic analyzes using dimensionality reduction. The method was fast, efficient, concise and able to analyze a large volume of data, not feasible for multiple sequence alignment. We created the first global phylogenetic tree of organisms carrying plastids using the plastome's coding regions, and we intend to make available both the tree and the analysis tool for scientific use.*

*Key-words: Phylogeny. Coding Region. SVect. Plastid.*

## LISTA DE FIGURAS

Figura 1: Tecnologias de sequenciamento.....	20
Figura 2: Início do processo de construção da árvore filogenética por NJ.....	26
Figura 3: Fluxograma do trabalho.....	28
Figura 4: Resultado da etapa de concatenação das proteínas.....	29
Figura 5: Esquema de funcionamento do algoritmo SVect.....	30
Figura 6: Gráfico de relação tempo/quantidade de sequências.....	34
Figura 7: Gráfico comparativo dos algoritmos para cálculo de distância.....	34
Figura 8: Comparação entre metodologia baseada em alinhamento e livre de alinhamento usando a matriz binária para obtenção de árvores filogenéticas.....	37
Figura 9: Comparação entre metodologia baseada em alinhamento e livre de alinhamento usando vetor de tamanho 800 para obtenção de árvores filogenéticas.....	38
Figura 10: Comparação entre metodologia baseada em alinhamento e livre de alinhamento usando vetor de tamanho 600 para obtenção de árvores filogenéticas.....	39
Figura 11: Posição da ordem Myrtales na árvore filogenética vetorial construída pelo método livre de alinhamento SVect.....	40
Figura 12: Árvore de BAYLY et al., (2013).....	41
Figura 13: Posição das ordens Rosales e Cucurbitales na árvore filogenética vetorial construída pelo método livre de alinhamento SVect.....	42
Figura 14: Análise filogenética de LI et al., (2016).....	43
Figura 15: Filogenia de ZHAO et al., (2017).....	44
Figura 16: Cladograma de RENNER; SCHAEFER, (2016).....	45
Figura 17: Posição da ordem Poales na árvore filogenética vetorial construída pelo método livre de alinhamento SVect.....	46
Figura 18: Árvore filogenética de SAARELA et al., (2018).....	47
Figura 19: Posição das ordens de Gymnosperms na árvore filogenética vetorial construída pelo método livre de alinhamento SVect.....	48
Figura 20: Parte da árvore filogenética de LU et al., (2014).....	49
Figura 21: Topologia do cladograma de YU et al., (2017).....	50
Figura 22: Posição das classes de algas vermelhas na árvore filogenética vetorial construída pelo método livre de alinhamento SVect.....	51
Figura 23: Reconstrução filogenética de YANG et al., (2015).....	52
Figura 24: Relações filogenéticas e rearranjo plastidial de PAIANO et al., (2018).....	53

## LISTA DE TABELAS

Tabela 1: Configuração de hardware dos computadores utilizados nos testes.....	31
Tabela 2: Tempo de processamento das abordagens de obtenção de árvores filogenéticas comparadas, rodando em um computador pessoal.....	33
Tabela 3: Tempo de geração de dendrogramas utilizando diferentes projeções em um hiper-espaço em dois computadores distintos.....	35

## LISTA DE QUADROS

Quadro 1: Proteínas usadas na análise controle.....	35
---	----



## LISTA DE SIGLAS

16S	<i>rRNA 16S ribosomal Ribonucleic Acid</i> (ácido ribonucleico ribossomal)
Aln	<i>Formato de arquivo de alinhamento</i>
DDBJ	<i>DNA database Japan</i> (Banco de dados de DNA do Japão)
DNA	<i>Desoxirribonucleic acid</i> (Ácido desoxirribonucleico)
ENA	<i>European Nucleotide Archive</i> (Arquivo de nucleotídeos Europeu)
faa	<i>Formato de arquivo multi-fasta de aminoácidos</i>
ftp	<i>File transfer protocol</i> (Protocolo de transferência de arquivos)
GenBank	<i>Genetic sequence database</i> (Banco de dados de sequência genética)
kpb	Kilo pares de base
ML	Maximum likelihood (Máxima verossimilhança)
MP	Maximum Parsimony (Máxima parcimônia)
NCBI	National Center for Biotechnology Information (Centro nacional para informação biotecnológica)
NGS	<i>Next generation sequencing</i> (Sequenciamento de próxima geração)
NIH	<i>National Institutes of Health</i> (Institutos nacionais de saúde)
NJ	Neighbour-joining (Junção de vizinhos)
OTU	<i>Operational Taxonomic Unit</i> (Unidade taxonômica operacional)
pb	Pares de base
RefSeq	<i>Reference Sequence Database</i>
SO	Sistema operacional
SSR	Simple Sequence Repeats (Sequência simples de repetição)
HTG	<i>Horizontal Transfer of Gene</i> (Transferência horizontal de genes)

## LISTA DE ABREVIATURAS

Fig.	Figura
Tab.	Tabela
pág.	Página

## LISTA DE SÍMBOLOS

®	Marca registrada
$\alpha$	Nível de significância
<	Menor que
=	Igual a

## SUMÁRIO

1.	INTRODUÇÃO.....	17
2.	OBJETIVOS.....	18
3.	JUSTIFICATIVA.....	18
4.	REVISÃO DE LITERATURA.....	19
4.1.	A Genômica .....	19
4.1.1.	Sequenciamento Genômico.....	19
4.1.2.	Montagem Genômica.....	20
4.1.3.	Anotação Genômica.....	21
4.2.	Os Plastídios .....	22
4.2.1.	O Cloroplasto .....	22
4.2.1.1.	O Genoma Cloroplastidial (cpDNA) .....	23
4.2.2.	Bancos de Dados para cpDNA .....	23
4.3.	Filogenética .....	24
4.3.1.	Métodos Filogenéticos .....	25
4.3.1.1.	UPGMA .....	25
4.3.1.2.	Junção de vizinhos ( <i>Neighbor-joining</i> ) .....	25
4.3.1.3.	Máxima Parcimônia .....	26
4.3.1.4.	Máxima Verossimilhança .....	26
4.3.1.5.	Inferência Bayesiana .....	27
4.4.	Inteligência Artificial .....	27
4.4.1.	Teorema de Johnson-Linnestras e Redução de Dimensionalidade .....	27
5.	MATERIAL E MÉTODOS.....	27
6.	RESULTADOS E DISCUSSÃO .....	32
6.1.	Desempenho .....	32
6.2.	Análise Filogenética Comparativa.....	36
6.3.	Filogenia Global de Plastídios .....	40
7.	CONCLUSÃO .....	53
8.	PERSPECTIVAS FUTURAS .....	54
	<b>REFERÊNCIAS .....</b>	<b>55</b>
	<b>APÊNDICE I: FILOGENIA GLOBAL DE PLASTÍDIOS COMPLETA .....</b>	<b>61</b>
	<b>ANEXO I: CERTIFICADO DE REGISTRO DE SOFTWARE SVECT.....</b>	<b>103</b>

## 1. INTRODUÇÃO

Com o advento do sequenciamento massivo paralelo de DNA e o consequente aumento exponencial de sequências nos bancos de dados devido ao baixo custo e redução de complexidade técnica, o campo da bioinformática tornou-se desafiador no que diz respeito a análises genômicas envolvendo *Big Data* (GREENE et al., 2014; EMMS; KELLY, 2015; YANG et al., 2017). Para STEPHENS et al., (2015), o campo da genômica será um dos maiores desafios da próxima década.

A filogenômica, campo que surgiu com expansão do volume de dados genômicos e proteômicos disponíveis, no ano 2000 (PHILIPPE et al., 2017) possui duas principais abordagens que utilizam os dados em larga escala para gerar reconstruções filogenéticas: super-árvores e super-matrizes (FLEISCHAUER; BÖCKER, 2017). Enquanto na abordagem de super-árvores são utilizadas sobreposições de árvores filogenéticas, em formato de matrizes que codificam os nós da árvore, para inferir a topologia final da grande árvore filogenética, a abordagem de super-matrizes vale-se de um grande número de genes ou proteínas dos organismos estudados para montar uma matriz que pode ser analisada pelos métodos convencionais de filogenia (FLEISCHAUER; BÖCKER, 2017), tais com o Máxima Verossimilhança – ML (FELSENSTEIN, 1981), Máxima Parsimônia – MP (FITCH, 1971) e Inferência Bayesiana (YANG; RANNALA, 1997). Há ainda uma terceira abordagem, conhecida como Mega-filogenia. Muito similar à abordagem de super-matriz, combina alinhamento de perfis com a seleção manual de regiões gênicas de interesse (SMITH et al., 2009).

Com o passar do tempo, reconstruções filogenéticas utilizando um grande conjunto de dados vem se tornando cada vez mais comum. Como exemplos, a primeira inferência filogenética em larga escala baseada no gene plastidial *rbcL* abrangendo 499 espécies de Viridiplantae (CHASE et al., 1993); o dendrograma que engloba 916 espécies de morcegos, publicado por JONES et al., (2002); o trabalho de MCMAHON & SANDERSON, (2006), que reconstruiu as relações filogenéticas entre 2.228 espécies de leguminosas, a árvore filogenética de 13.533 espécies de plantas verdes baseadas no gene *rbcL* de SMITH et al., (2009); PYRON et al., (2013) ao gerar um dendrograma para Squamata com cerca de 4100 espécies de cobras e lagartos baseado em um conjunto de 12 genes; FITZJOHN et al., (2014), que reconstituiu as relações filogenética de aproximadamente 24 mil espécies de plantas lenhosas; a análise de 37 genomas completos para inferir a filogenia de Oomicetos publicado por MCCARTHY & FITZPATRICK, (2017); e a grande árvore de Viridiplantae de SMITH & BROWN, (2018), contendo 79.881 espécies.

Em seu trabalho, DE PIERRI (2017) propôs um algoritmo livre de alinhamento para comparação de sequências que foi aplicado em uma reconstrução filogenética para 6.811 proteomas mitocondriais, o SVect.

## 2. OBJETIVOS

### 2.1. Objetivo Geral

Propor e validar a aplicação da representação vetorial de sequências proteicas e a projeção em um hiperespaço como estratégia para identificar grupos de sequências que compartilham algum grau de semelhança em sua composição de aminoácidos, como uma eficiente alternativa para análises filogenéticas em grandes bases de dados, tais como banco de genomas completos, onde o alinhamento múltiplo de sequência é inviável.

### 2.2. Objetivos Específicos

Entender o desempenho computacional do algoritmo SVect e encontrar o melhor conjunto de parâmetros para análises filogenéticas baseadas em proteomas de plastídios;

Comparar alinhamento de sequências com o SVect;

Empregar o SVect no conjunto de sequências codificantes traduzidas dos portadores de genoma plastidial, a fim de obter uma filogenia global.

Validar o dendrograma global gerado utilizando o SVect na literatura;

## 3. JUSTIFICATIVA

Para análises filogenéticas em geral, a limitação dos métodos baseados em alinhamento múltiplo para a comparação de sequências se torna evidente devido à complexidade algorítmica para se obter um alinhamento ideal (VINGA; ALMEIDA, 2003). Para tornar o método viável, poucas sequências de aminoácidos ou nucleotídeos são utilizadas para um grande número de organismos (vide introdução) ou poucos organismos são necessários para se utilizar de genomas completos em análises. Podemos perder informação evolutiva ao usarmos um número mínimo de genes/proteínas/sítios ao inferir relações filogenéticas (DE PIERRI, 2017), ou o prejuízo pode ser sofrido no tempo computacional para obter os resultados da análise. Mediante a tais limitações da abordagem utilizando o alinhamento múltiplo de sequências, é evidenciado a

necessidade de um método capaz de gerar filogenias consistentes a partir de grandes volumes de dados em tempo computacional não-proibitivo.

## 4. REVISÃO DE LITERATURA

### 4.1. A Genômica

Por meio de avanços científicos, a genômica tem propiciado cada vez mais investigações a nível molecular (SHENDURE et al., 2004). A partir do momento em que o DNA foi descoberto e descrito na década de 1950, muitas tecnologias foram criadas e a sociedade em geral tem ganhado com esses avanços, que englobam entre outras áreas, a economia e a medicina (WATSON; CRICK, 1953; DINIZ, 2007).

A criação do projeto genoma humano em 1990 foi um símbolo de conquista de conhecimento no contexto da genômica, que passou a ser uma ciência relevante na área entre o meio das ciências biológicas (CATANHO et al., 2007; DINIZ, 2007). O emprego dos genomas completos ou parciais auxiliam no processo de comparação para entender a conformação de outro genoma (SHENDURE et al., 2004).

A Bioinformática aplicada à genômica é uma ciência que utiliza e maneja os dados oriundos do sequenciamento de amostras genéticas de diversas espécies de organismos. Isto posto, os dados resultantes do método de sequenciamento que estão fragmentados e embaralhados são arranjados afim de dar um entendimento mínimo em relação à configuração genômica de um determinado organismo (SIM et al., 2012).

Nesta sessão serão rapidamente revistos os passos para a obtenção de um genoma completo, um dos principais objetivos da genômica.

#### 4.1.1. SEQUÊNCIAMENTO GENÔMICO

A etapa inicial para que possamos realizar a análise do genoma é o sequenciamento. A técnica de sequenciamento de genomas que consistia em fragmentar o DNA em pequenos pedaços conforme o tamanho do vetor usado para a clonagem era conhecida como técnica de Sanger. Essa técnica consiste em alongar primers de RNA sob uma fita de DNA utilizando didesoxinucleosídeos trifosfato além dos 4 nucleotídeos comuns. Assim, quando um didesoxinucleosídeo era adicionado na molécula em construção, interrompia a síntese (Fig 1A). Então, as bases que interrompiam a biossíntese eram identificadas por meio de eletroforese.

Essa técnica tornou-se possível após avanços tecnológicos da metodologia shotgun. (SANGER; COULSON, 1975)

Com as técnicas de sequenciamento de próxima geração (NGS), sequenciamento de segunda geração, ou sequenciamento massivo paralelo, houve uma modernização dos processos a ponto de dispor de uma clonagem *in vitro*, dispensando a clonagem laboratorial. Estes avanços, como a multiplexação e a marcação fluorescente de nucleotídeos (Fig 1B) ao invés dos métodos eletroforéticos resultaram em menores custos, maior velocidade de obtenção das sequências e um ganho considerável na qualidade do sequenciamento, além de chegar a sequenciar bilhões de pares de base em uma corrida (CARVALHO; SILVA, 2010; SHENDURE et al., 2017).

Atualmente, o sequenciamento em tempo real, de terceira geração ou de molécula única (diagramado simploriamente na Fig. 1C) reduziu ainda mais os custos e a facilidade de produção de sequências, visto que as amostras podem ser preparadas em campo. Além disso, um salto na qualidade das sequências foi conquistado, já que o tamanho dos fragmentos aumentou de cerca de 300-400pb nos sequenciadores de segunda geração, para aproximadamente 150kpb nos sequenciadores de terceira geração (BLEIDORN, 2016).

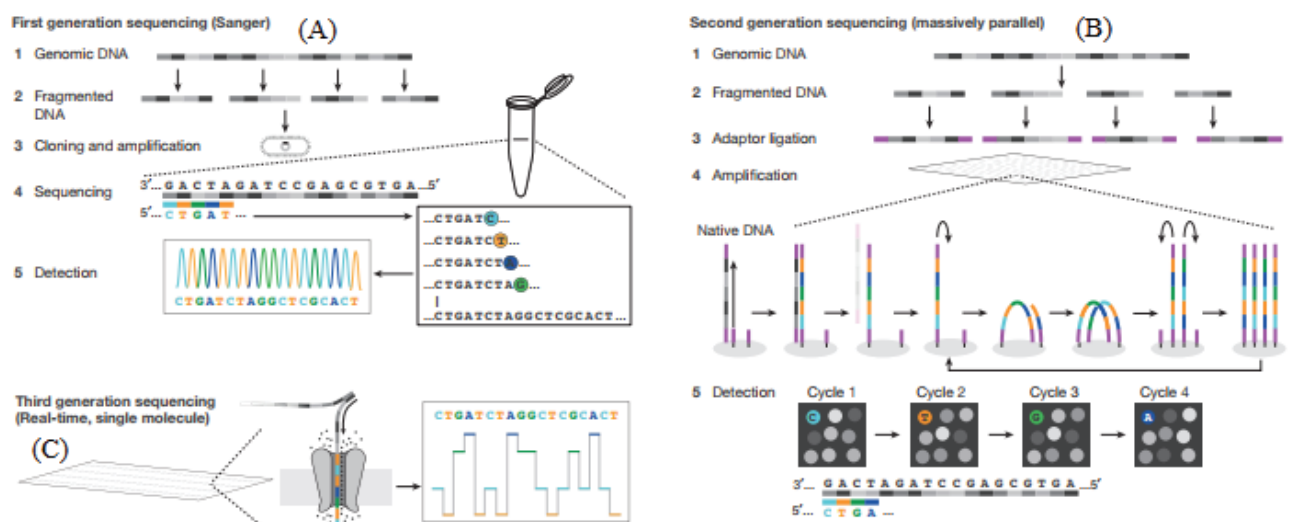


Figura 1: Tecnologias de sequenciamento. O esquema exemplifica o funcionamento dos sequenciadores de primeira, segunda e terceira gerações. Fonte: (BLEIDORN, 2016).

#### 4.1.2. MONTAGEM GENÔMICA

A montagem de genomas é realizada por meio de agrupamento, em que sequências do DNA são compiladas utilizando como critério a sobreposição das *reads*. As *reads* são definidas



como fragmentos de sequência que são originados pela quebra do DNA no sequenciamento. São utilizados programas de Bioinformática para realizar o agrupamento das *reads* por superposição de sequências idênticas, a fim de reconstruir o fragmento original (SIMPSON; POP, 2015).

Existem hoje muitos programas que realizam a montagem genômica, utilizando diferentes tipos de algoritmos para obter a ordem das sequências do DNA. Mesmo com uma gama de programas, ainda existem problemas que podem interferir negativamente na montagem, como agrupamentos repetidos, regiões com baixa qualidade de sequenciamento, o tamanho dos fragmentos e o próprio montador em si (SIMPSON; POP, 2015).

Os procedimentos que se baseiam no percentual de identidade para o reconhecimento de sequências de outros genomas para realizar a montagem são confiáveis, se o material norteador for de qualidade. Os elementos obtidos através da anotação gerados a fim de serem usados para mineração dos dados em banco de dados públicos, confrontando com genes já existentes, concluindo a montagem de forma adequada. É importante lembrar que a utilização de ferramentas de bioinformática deve ser realizada de forma integrada, para que a verificação de dados importantes seja favorecida (TROY et al., 2001).

#### 4.1.3. ANOTAÇÃO GENÔMICA

Após a realização do sequenciamento, o próximo passo para a análise do genoma consiste na anotação e a busca por similaridade constitui o início dessa etapa. Neste ponto do trabalho, são usadas as sequências homólogas a fim de treinar os métodos estatísticos, para que esses sejam capazes de realizar o reconhecimento de genes. Isso resulta em um arquétipo que será empregado para analisar o restante do genoma (TROY et al., 2001).

Para STEIN (2001), há três categorias na anotação de genoma que se destacam por serem mais relevantes: a anotação em nível de nucleotídeos, ou seja, quais os nucleotídeos que compõem a sequência; anotação em nível de proteína, que determina quais proteínas são codificadas pelos genes; e anotação em nível de processo, definindo que função a proteína desempenhará no organismo. Esses processos devem ser realizados de maneira a alcançar resultados de alta qualidade, para que se possa identificar corretamente os genes, os produtos eles codificam, ratificando as características em particular.

Montagens erradas podem levar à identificações erradas de genes e a anotações equivocadas, e se levarmos em consideração o fato de que futuras pesquisas também correm o risco de serem afetadas pelo erro cometido, a dimensão do problema pode fugir ao alcance. Mesmo

utilizando um bom preditor de genes e seguindo um controle de qualidade, ainda assim a anotação está sujeita a erros (YANDELL; ENCE, 2012) .

#### 4.2. Os Plastídios

Os plastídios são organelas membranosas que estão presentes em diversos grupos de organismos eucarióticos como as plantas, protozoários e algas. Estas organelas são classificadas conforme a função que exercem no organismo. São elas os proplastídios que são encontrados nos meristemas da planta e dão origem aos outros plastídios; como etioplastos, cloroplastos e leucoplastos. Os etioplastos são plastídios não portadores de clorofila, encontrados em cotilédones de plântulas de angiospermas e que podem ser convertidos em cloroplastos; Os cloroplastos são plastídios portadores de clorofila e realizam a fotossíntese; os cromoplastos são plastídios coloridos, derivados dos cloroplastos e funcionam como depósitos de carotenoides e outros pigmentos, estão presentes em pétalas de flores e frutas; Amiloplastos são derivados do leucoplastos e encontrados em tecidos de armazenamento, como tubérculos e endospermas de sementes; os elaioplastos são também derivados dos leucoplastos e cuja função é armazenar lipídeos; e os leucoplastos são plastídios sem pigmento encontrados nas células da raiz e dão origem aos amiloplastos, elaioplastos e proteinoplastos , ( LOPEZ-JUEZ & PYKE, 2005 e EGEA et al., 2010).

Assim como as mitocôndrias, acredita-se que os plastídios foram resultado de um evento endossimbiótico. Um procarionte gram-negativo fotossintetizante foi fagocitado por um eucarioto anaeróbico, e a bactéria não foi digerida passando a trabalhar em conjunto com o hospedeiro. Acredita-se que parte do genoma do procarionte foi perdido ao longo da evolução. Outra parte desses genes foi incorporada ao genoma nuclear da célula eucariótica, tornando a cooperação irreversível e indispensável para a sobrevivência (WEEDEN, 1981; GOULD et al., 2008).

##### 4.2.1. O CLOROPLASTO

Cloroplasto é uma organela de formato variado, revestido por dupla membrana. No estroma, plasma interno análogo à matriz mitocondrial, se localizam todos os orgânulos e enzimas metabólicas (BLANKENSHIP, 2013). Uma das estruturas presentes no estroma é chamado de tilacóide, que geralmente estão empilhados em forma de grana. O tilacóide é revestido por uma terceira membrana. Através desta membrana, um processo semelhante à cadeia transportadora de elétrons mitocondrial gera ATP e NADPH, que são usados como energia para realizar o metabolismo da planta (ANDERSON, 1975).

Além da produção de energia celular no tilacóide, uma enzima presente no estroma chamada *ribulose-1,5-bifosfato-carboxilase*, em plantas processa o gás carbônico atmosférico para produzir açúcares na fase escura da fotossíntese, que serão modificados posteriormente no citoplasma, servindo como fonte de energia para outras partes da planta e como substância de reserva (CLELAND et al., 1998).

#### 4.2.1.1 O GENOMA CLOROPLASTIDIAL (CPDNA)

O cloroplasto, de maneira similar à mitocôndria, também possui o seu próprio material genético (BURGESS, 1985). Em um cloroplasto jovem, podemos encontrar até 100 cópias do cpDNA no estroma, porém, de acordo com a idade, o número de cópias pode ser reduzida a 15 (HELDT; HELDT, 2005). É uma dupla fita comumente circular, tendo em média 150.000pb (CLEGG et al., 1994), podendo variar entre 120.000 e 170.000pb (SHAW et al., 2007) e contém *introns* (KREBS et al., 2014). Possui uma média de 120-130 genes (ROGALSKI et al., 2015; DANIELL et al., 2016) sendo o menor cpDNA composto de 4 genes (*Monoraphidium neglectum*) e o maior, de aproximadamente 273 - *Pinus koraiensis* (NCBI, 2018).

O cpDNA de plantas terrestres possui uma estrutura muito semelhante. Eles codificam 4 rRNAs, 30-31 tRNAs, essenciais para a síntese proteica dos cloroplastos, 21 proteínas que se associam ao ribossomo e 4 subunidades da RNA polimerase, além de 28 proteínas tilacóidais com função conhecida no sistema fotossintético (como a ATP sintase e a NADH desidrogenase) e a ribulose-1,5-bisfosfato-carboxilase, principal enzima do sistema (KREBS et al., 2014).

#### 4.2.2 Bancos de Dados para cpDNA

O GenBank, vinculado ao NCBI (BENSON et al., 2017) é um banco de dados biológico de sequências genômicas mantido pelo NIH, e que faz parte de uma cooperação internacional com o DDBJ e o ENA. Essas três bases de dados trocam informações e são mutuamente atualizados todos os dias. O Genbank possui dados de todas as esferas de vida, sendo possível encontrar cpDNA e proteínas plastidiais, porém, os dados são redundantes e em sua maioria não passam por curadoria.

O ChloroplastDB (CUI et al., 2006) é uma base de dados especializada em cpDNA. Atualizado e mantido pelo GenBank, ele aplica técnicas de mineração genômica. A proposta do ChloroplastDB é disponibilizar ferramentas com uma interface amigável que facilite a análise dos dados de sequências cloroplastidiais, entretanto, esse banco está desatualizado, com sua última atualização em 2007, cujos números não ultrapassavam 100 genomas completos.

O GObase (O'BRIEN et al., 2009) é um banco de dados para genomas organelares que engloba plastídios e mitocôndrias e foi criado em 2008. Também se alimenta de dados disponibilizados pelo GenBank e possui diversas ferramentas para manipulação, análise e mineração de dados. Todavia, O GObase também se encontra desatualizado, com sua última atualização datando de 2010.

Outros dois bancos de dados importantes são o ChloroMitoSSRDB 2.0 (SABLOK et al., 2013, 2015) e o ChloroSSRdb (SHANKER, 2014). Com ênfase em microssatélites e SSRs, ambos contemplam genomas completos, no entanto, tais repositórios contêm anotações apenas de SSRs e microssatélites, não fornecendo anotações completas de genes (proteínas e demais produtos).

O *Plastid-encoded Protein Clusters* (ZVERKOV et al., 2012, 2015; LYUBETSKY et al., 2013) é um banco de dados de clusters proteicos de plastídios que aplica mineração em todos os organismos que contêm plastídios sequenciados. O banco está sendo constantemente atualizado, entretanto, não disponibiliza os dados para download. Como o *Plastid-encoded Protein Clusters* é apenas para consulta, o pesquisador não consegue expandir suas pesquisas, ficando limitado às ferramentas do banco.

O Organelle Genome Resources (WOLFSBERG et al., 2001), disponibiliza para descarga os dados do RefSeq (PRUITT et al., 2002, 2007; O'LEARY et al., 2016) – genomas, RNA's codificados e proteínas traduzidas, tanto para plastídios, quanto para as mitocôndrias. Esse banco é bimestralmente atualizado pelo NCBI e por esse motivo escolhemos trabalhar com essa base.

### 4.3 Filogenética

A filogenética, também conhecida como cladística, é a ciência que tem por objetivo criar classificações biológicas baseadas primordialmente em sequências genéticas, inferindo relações de parentesco e ancestralidade entre os organismos ao longo do tempo. O resultado de uma análise filogenética é um cladograma, que é uma representação hierárquica das relações evolutivas em forma de uma árvore binária. Dessa forma, um ramo da árvore representa um grupo de organismos que tem um mesmo ancestral comum, sendo denominado clado ou grupo monofilético (PATWARDHAN et al., 2014; NICOLAU, 2017).

Há diversos métodos para se obter uma árvore filogenética, e os principais estão sendo tratados de maneira resumida no subtópico a seguir.

#### 4.3.1 Métodos Filogenéticos

Os métodos filogenéticos são genericamente divididos em 2 grupos: os baseados em matrizes de distâncias e os baseados em características, ou dados discretos (PATWARDHAN et al., 2014). Os métodos baseados em distância tomam como premissa que a distância evolutiva real entre duas espécies pode ser calculada pela diferença entre suas sequências. Esses métodos empregam o número de alterações entre pares em um grupo de sequências para produzir uma árvore filogenética do grupo (FELSENSTEIN, 1996), e são subdivididos em baseados em algoritmos de agrupamento, e baseados em otimização (PATWARDHAN et al., 2014).

#### 4.3.1.1 UPGMA

Um dos métodos baseados em matrizes de distâncias, o UPGMA, do inglês, *Unweighted Pair Group Method with Arithmetic mean* ou método não ponderado de agrupamento por média aritmética, consiste em um algoritmo de agrupamento (*clusterization*) hierárquica de baixo para cima (*bottom-up*), ou seja, partindo das OTU's visa unir os grupos mais próximos criando um novo agrupamento até chegar ao todo, unindo o conjunto em um único grande grupo (SOKAL, R; MICHENER, 1958). É o algoritmo mais simples, e gera arvores enraizadas ultramétricas que refletem a estrutura da matriz de entrada.

#### 4.3.1.2 Junção de vizinhos (*Neighbor-joining*)

Este método de junção de vizinhos, é um algoritmo de agrupamento de baixo para cima *que* considera como a melhor árvore, aquela que tem a menor soma total dos ramos (SAITOU; NEI, 1987). O NJ considera cada organismo como portador de uma taxa de evolução distinta, excluindo a hipótese do relógio molecular (ZUCKERKANDL; PAULING, 1965). Essa é uma vantagem desse método em relação ao UPGMA. O NJ consiste em 5 passos, segundo seus criadores:

- 1- Com base na matriz de distância atual, calcule a matriz Q.
- 2- Encontre os dois tipos distintos de i e j (ou seja, com  $i \neq j$  para qual Q (i, j) tem o menor valor possível. Esses organismos estão unidos a um nó recém-criado, que é conectado ao nó central. Na Figura 2-B, as OTU's são unidos ao novo nó X.
- 3- Calcule a distância de cada um dos organismos do par para este novo nó.
- 4- Calcule a distância de cada um dos organismos fora desse par para o novo nó.
- 5- Comece o algoritmo novamente, substituindo o par de vizinhos unidos com o novo nó e usando as distâncias calculadas na etapa anterior.

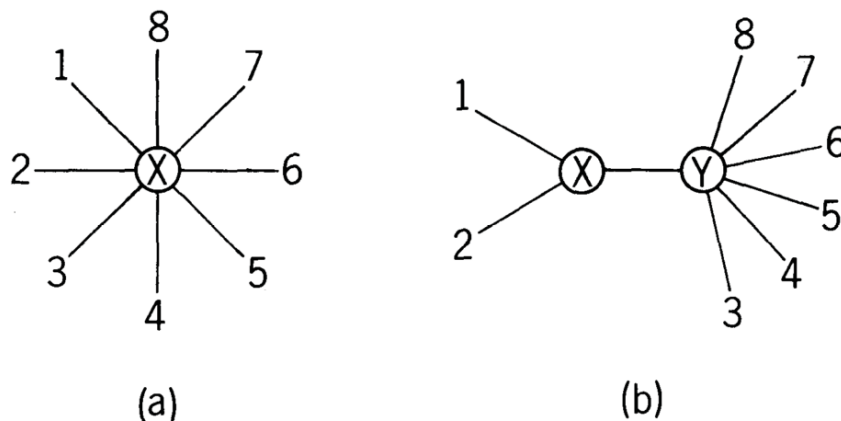


Figura 2: Início do processo de construção da árvore filogenética por NJ. (A) Árvore em formato de estrela sem hierarquia. (b) OTU's 1 e 2 unidas pelo nó X após o primeiro ciclo. Fonte: Adaptado de Saitou & Nei (1987).

Já os métodos baseados em características ou dados discretos, também conhecidos como métodos de busca, examinam individualmente cada coluna do alinhamento, procurando a árvore que melhor represente toda a informação. Destacamos aqui o Máxima parcimônia (DAY, 1987), Máxima verossimilhança (FELSENSTEIN, 1981) e a Inferência bayesiana (RANNALA; YANG, 1996). O funcionamento desses algoritmos são temas dos subtópicos a seguir.

#### 4.3.1.3 Máxima Parcimônia

Máxima parcimônia é um método estatístico não-paramétrico que minimiza o número de alterações em uma árvore filogenética conferindo situações de caracteres a nós internos na árvore. O comprimento do caractere (ou do sítio) é o número mínimo de mudanças necessárias para esse sítio, enquanto a pontuação total da árvore é a soma dos comprimentos de caracteres em todos os sítios. A árvore de máxima parcimônia é a árvore que torna mínima a pontuação da árvore, ou seja, evolutivamente falando, a melhor árvore é aquela que exige o menor número de mudanças para explicar os dados observados (DAY, 1987; SWOFFORD, 2002). Uma ferramenta de destaque na utilização dessa abordagem é o PAUP (SWOFFORD, 2002).

#### 4.3.1.4 Máxima Verossimilhança

Em uma análise filogenética por máxima verossimilhança, percebemos duas etapas de otimização: a otimização dos comprimentos dos ramos para calcular a pontuação da árvore e a otimização da busca pela árvore de máxima verossimilhança. Como é um método estatístico, a topologia da árvore torna-se um modelo, enquanto os comprimentos dos ramos e os parâmetros de substituição compõem os parâmetros desse modelo. Esse método consiste em comparar diversos modelos estatísticos com a mesma quantidade de parâmetros para encontrar a

topologia de máxima verossimilhança (FELSENSTEIN, 1981; YANG, 1996). Um software popular que implementa essa metodologia é o RAxML (STAMATAKIS, 2006).

#### 4.3.1.5 Inferência Bayesiana

A inferência bayesiana é uma metodologia geral de inferência estatística. Quando comparada à máxima verossimilhança, os parâmetros no modelo são considerados variáveis aleatórias com distribuições estatísticas, diferindo do ML, que possui parâmetros constantes, fixos e desconhecidos. Antes de analisar os dados brutos, estes são combinados a uma distribuição anterior, para então estimar a distribuição posterior (RANNALA; YANG, 1996). Uma das aplicações mais populares da inferência bayesiana foi a implementação do algoritmo MCMC, *Markov chain Monte Carlo* ou Cadeia de Markov de Monte Carlo (YANG; RANNALA, 1997), implementado no software MrBayes (HUELSENBECK; RONQUIST, 2001).

O Mega 7.0 (KUMAR et al., 2016) é um pacote amplo, que disponibiliza tanto os métodos baseados em distância, quanto os métodos baseados em características.

### 4.4 INTELIGÊNCIA ARTIFICIAL

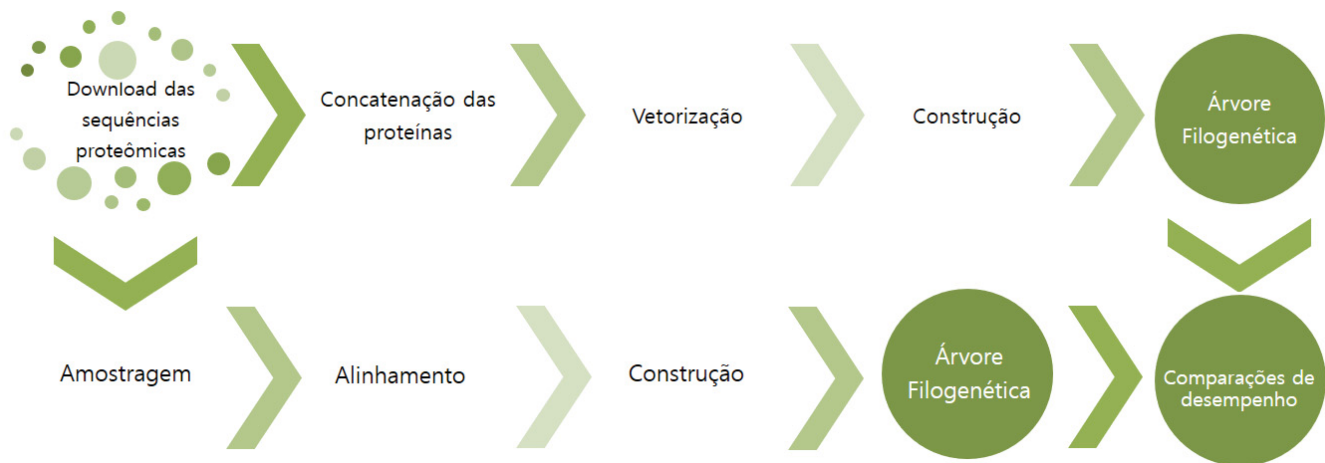
#### 4.4.1 Teorema de Johnson-Linnestraus e Redução de Dimensionalidade

Esse teorema afirma que um pequeno conjunto de dados significativos em um espaço de alta dimensão pode ser projetado em um espaço consideravelmente menor, de forma que as distancias proporcionais entre os pontos é praticamente preservada. É possível manter a qualidade de um grande volume de dados, ao multiplicá-lo por um conjunto de vetores ortonormais aleatórios. Isso reduz o volume de dados a ser trabalhado computacionalmente, sem perder as propriedades dos dados originais (JOHNSON; LINDENSTRAUSS, 1984; FRANKL; MAEHARA, 1988)

## 5 MATERIAL E MÉTODOS

O trabalho foi desenvolvido de acordo com o fluxograma da Figura 3. Cada subtítulo descreve uma etapa de trabalho representado na imagem.





**Figura 3: fluxograma do trabalho. O trabalho foi desenvolvido sob duas perspectivas: a construção da grande árvore a partir da aplicação do método proposto e a avaliação da consistência do resultado final, bem como uma simples análise de desempenho ao final do processo. Fonte: O autor.**

### 5.1. Obtenção das Sequências Proteicas

Inicialmente, obtivemos as sequências proteicas de todos os organismos com proteoma cloroplastidial disponível nos arquivos do RefSeq (PRUITT et al., 2002, 2007; O'LEARY et al., 2016) por meio de ftp, acessados pelo Organelle Genome Resources (WOLFSBERG et al., 2001) no formato \*.FAA (Fasta de aminoácidos). A captura dos dados era realizada sempre que uma versão atualizada do banco era lançada, bimestralmente.

### 5.2 Concatenação das Proteínas

Em posse dos dados, concatenamos então as sequências de proteínas pertencentes ao mesmo organismo, inserindo um separador de cinco asteriscos (\*\*\*\*\*) entre uma proteína e outra, fato que será explicado no próximo parágrafo. Realizamos ainda a enumeração dos organismos por meio de um identificador sequencial para evitar problemas com dados duplicados, visto que os cabeçalhos das proteínas foram descartados. O resultado dessa etapa é exemplificado na Figura 4 :



```

>Chromera velia org 2222
MLFQRCTPLYLYKYSNTRLGSGKKRQPLSVSRNLPKSKPARSYEVLYFIALVLLGCGLLYVLLSPLQKLGL
ARTSKRFFASLNLDAKTAKNVINCLVRIKSIARRGEVTLIEDMLVIYKHIDHTLDVLSSPRHLVSLLPYDF
RALEVENIAGVLGVGPSAAKSLAEIFLTPVGNTTSSALVFQSEWSGLLPFGADYLMCLNAERQQVAFWWT
LSALKEIGPARQWDFTPRHTLGIPKKLAELKYKYLEEAAEALGQLNYIADQWQ*****MKTKRFTWMRE
ISRTVFNFVQNTDRRLYIGWFGLLMVPTLLAVSCFIAAFVAAPPVDIDGIREPVSGSLLYGNNIISGAV
VPSSNAIGMHLYTVWDAMCIEEWLYNGGPYQFVVLHFLIGVASYLGREWELSYRLGMRPWFVAFSAPVA
AASAVFLVYPIGQGSFSDGMPLGISGTFFNMLVFQAEHNILMHPFHMAGVAGVFGGSLFSAMHGSVLTST
LISMTSEDESPNYGKFGQYETYNIVAAHGYFGRLLIFQFASFNNRSRLHFFLAIWVVGIVLTAMGIST
MAFNLNGFNFNQSIIDAEGRVISSWADILNRANLGMEVMHERNAHNFPLDLA*****MEGIVAACSAFSA
GIAIGLGSIGPGIGQILAGDAVSAISRQLEGEDKIRTLLPSLAVLEAVTIYGLLIALVIGRVALLKVD
DSSKNPDQNSHEEKAIEPCYVVGAACSAFSYASVIDIIDVLCPLNRTKNKEPKRDLVWYLQLIWYLIRY
LYFSAYWKYSNVNREW*****MLLSTFSRYTRRPLHKHAPPFRILNTTASILETGKVIDLVTPIYRIGGK
IGLFGGAGVGKTVVIMELITNIAKLHGKISVFAGVGRSREGFDLYKEMRESGVIDTVFPDPSKVSIVYG
QMNETPAARMRVALGAIVAEAFRDEFFQNVLFIDNIFRYVQAGAEISALLGRIPSAVGYQPTLATMG
SLQERIASVSRLSLESQMPIISSITSIQAIYVPADDLTPAPATNFAHLDATTVLSRNLAQKGIYPAIDP
LASSSLCLKPVECGPFHYSITQKIIGVLSARNIQDLLSIIGFEELSEEQKITIYRSRKIQNFLSQPFFV
AEVFNNIPGKLVNLSETVIDAQLILNGYADAIEDISLSYVGRLEAIQKDYKNFYKFFIVDKELFLTGW
LKEKAQIVDETFYRNDPQTKAPRNLVEDLDQKVLNEDGSLPPLSRIWWLAERPGLVTNHLNLSVEDLLEE
LALPPESRTLPLSL*****MIITTTIIAPKSVPLSSVGKNWESWKSNDNGIRFLERKVVEMLLITFSRYT
RRRLNKHAPNILFSRRTPIISLQARQQSGMKNLAQQERVAVSLPQKPATRTGSSEILGLVALFMVVSFVI
AILYVRFSAGPKVGFAGNSKVVGVDLPWVHRKRKFLGEIIGRFVAAMANWARDCKMYRRYIKYGWYP
PPSPRLDEILYQFVELWEFCRQVYQDYLRSKTEGRPPV*****MDIEDPVQTFGEYEFDLDEVVYRLQH
LLFLSVQILCLSFIVSPIVSDVVFETWYFKEVSFMGQSVNTIFFNFFVSFLLSVVILSSYIVFQIFLWL
DTAFLEDERLLYRFLNGFWLFFFLALGAAIFLIFPLYWITVLQGIKFLNIQNFQSIVDLSQFTDLAFL
VFSTWWSLTPILQIVILYNLITVENLISNWKYVIIIFSMLGCGILTASADPLQLLALVINLFFIFSL
ILYKSLFHNMCVLFDEKSLTTEDDDDLFLMFLWE*****MFANFPPIYAQQGNLIAREANGKLVCANCLN
NYPKIYGVRLHNLNEVFNKLKSVDRPNSALTQSASLNGSKKFLGLGGIVILPEQFDVYRPAENPFIP
YSTEEDQTLVVGPLYTNKATLNIIPVRSPAAGDGIPTTNMFFGGINRGRGQLTPLGVKTDITNKPELEFSNL
SVGQIWRSSKGHLYMTQANSENLYRLNGLHLTWGDNNAVIVNSPNQGGFGQSEFSTSVLNLEYLVNLYLF
FTLFCVSAQISLIIYKKDYLVRSSFDTWSYWCRAFALPKRKIVVSEETVMVDPNFIVLTGKKYKR

```

Figura 4: Resultado da etapa de concatenação das proteínas. Um arquivo multifasta é gerado com a seguinte configuração: Cada cabeçalho contém o nome da espécie a qual pertencem as sequências proteicas (*Chromera velia*) e um identificador sequencial (org 2222), enquanto as sequências se encontram concatenadas por cinco asteriscos (\*\*\*\*\*). Nenhuma ordenação foi usada na concatenação, sendo que a primeira proteína do concatenado é a primeira proteína do arquivo pertencente ao organismo em questão, e assim sucessivamente. Fonte: O autor.

### 5.3 Vetorização e Redução de Dimensionalidade

Em seguida, utilizamos o algoritmo SVect (DE PIERRI, 2017) para criar uma representação vetorial de cada proteoma. Esse algoritmo usa uma janela de tamanho cinco ( $k=5$ ), com a terceira posição vazia, dividindo os aminoácidos selecionados em dois pares de dois resíduos. A janela percorre a sequência, conforme a Figura 5 (A) e uma matriz é gerada, de acordo com a ocorrência desses pares relacionados (Figura 5 B).

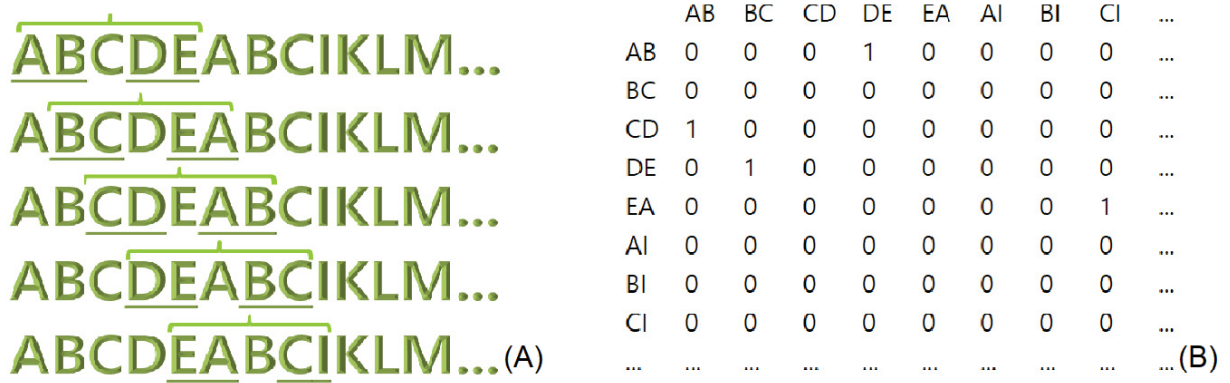


Figura 5: Esquema de funcionamento do algoritmo SVect (ir para a pág. 102). Uma janela móvel (k-mer) de tamanho cinco e espaçada na terceira posição percorre toda a sequência (A), marcando todas as combinações de tetra-aminoácidos em uma matriz binária quadrada de tamanho quatrocentos. Como resultados, temos uma matriz  $M_{m \times n} = M_{400 \times 400}$  (B), tal que  $m$  e  $n$  são todas as possíveis combinações para pares de aminoácidos tradicionais. Os espaçadores entre as sequências de resíduos de aminoácidos inseridos na etapa anterior foram usados como recurso para evitar a criação de dados não-biológicos na matriz resultante. Fonte: O Autor.

Concatenamos então as colunas da matriz para gerar um vetor de 160 mil posições binárias. Essa quantidade de memória pode ser um problema quando se trata de um grande volume de dados. Implementamos então um algoritmo simples que aplica o teorema de Johnson-Linnestraus para reduzir o tamanho dos dados em memória. Após alguns testes, para fins de divulgação, usamos o vetor binário sem redução de dimensão e a projeção em um hiper-espço de tamanho 600 para cada proteoma. Os testes que foram norteadores para a tomada de decisão encontram-se na sessão de Resultados e Discussão.

#### 5.4 Construção

Com os proteomas vetorizados e projetados, extraímos a matriz de distâncias internas entre os pontos, usando o cálculo de distância euclidiana. Uma implementação do método de junção de vizinhos foi usada para a obtenção da árvore filogenética, dada a matriz de distancias como entrada para o algoritmo. Todos os passos realizados até então se desenvolveram em ambiente *MatLab*®, utilizando os parâmetros padrões.

#### 5.5 Amostragem

Já em posse da filogenia, utilizamos scripts em Python3® para eleger 26 genes presentes em 21 organismos no conjunto de dados brutos não processados oriundos do RefSeq. Essas proteínas estão listadas no quad. 1. Ordenamos e concatenamos as proteínas selecionadas por espécie, utilizando os mesmos espaçadores da primeira etapa desse processo (\*\*\*\*\*). Todas as sequências foram salvas em um multi-fasta de aminoácidos.

5.6 Alinhamento Múltiplo do Grupo Controle

A etapa anterior nos forneceu dados bem controlados, afim de comparar a eficiência do método livre de alinhamento com o alinhamento de sequências tradicional. O alinhamento foi realizado em SO Linux, utilizando o algoritmo *Clustal Omega*® usando 20 *threads*, e o resultado foi direcionado para um arquivo \*.aln.

5.7 Obtenção da Árvore Controle

Nosso arquivo de alinhamento foi então carregado no programa *SeaView4*® para que a matriz de distâncias fosse gerada a partir do cálculo de *Poisson*. A árvore filogenética também foi gerada no *SeaView4*®, usando o algoritmo junção de vizinhos (*neighbor-joining*). Foram rodados 1500 ciclos de reamostragem aleatória (*bootstrap*) para avaliar a consistência da árvore, e o consenso foi usado para as comparações.

O procedimento vetorial do item 3.2.3, também foi usado nesse conjunto amostral para obter um dendrograma. E de forma comparativa, pudemos avaliar a confiabilidade do método utilizado no trabalho.

5.8 Desempenho

Realizamos também testes de desempenho para demonstrar a confiabilidade da estrutura das árvores filogenéticas obtidas pelo algoritmo *SVect*. Inicialmente, aferimos o tempo de processamento para a geração dos dendrogramas construídos a partir do *Svect*, baseados nas matrizes binária e projeções em um hiper-espaco de 1200, 1000, 800, 600, e 400 atributos. Pare esse teste, foram usados os parâmetros padrões do *MATLAB*® para a função *NJ* e demais.

Então, aferimos o tempo de processamento para, a partir da sequência, construir uma arvore filogenética, incluindo todos os pré-processamentos necessários. Dessa forma, comparamos as metodologias baseadas em alinhamento e livre de alinhamento no requisito velocidade de processamento. Dois computadores foram usados para este fim, e a descrição de hardware se encontra na Tabela 1.

Tabela 1: Configuração de hardware dos computadores utilizados nos testes

Servidor	Notebook
Processador Intel® Xeon® E5-2640 v4 40- Core @ 2,4GHz 256 GB RAM	Processador Intel® Core I7-4510U ® 2.6GHz 8 GB RAM

2 TB HDD

1 TB HDD

Linux 4.4.0-72 #93 Ubuntu SMP

Windows® 10 Pro

Fonte: O Autor.

Determinamos o tempo de processamento computacional a medida em que o número de de organismos com regiões codificantes inseridas na análise aumenta, visando entender a dinâmica de processamento do algoritmo. Geramos árvores filogenéticas com o mesmo conjunto de parâmetros para dados de proteomas liberados pelo RefSeq no período de julho de 2016 até dezembro de 2017. Por fim, buscamos entender a relação entre a matriz binária e as projeções em um hiper-espço, para fixar um ponto de equilíbrio entre desempenho e qualidade da árvore. Geramos 30 projeções em um hiper-espço entre 100 e 3000 atributos, incrementando de 100 em 100 para os quatro principais algoritmos de cálculos de distância que são aplicáveis a esse caso: distância euclidiana, distância de cossenos, distância de correlação e distância de Spearman. Utilizamos então a correlação de Pearson entre a projeção em um hiper-espço peculiar em tamanho e tipo de distância usada (por exemplo, tamanho 1000 gerado a partir de distância dos cossenos) e a matriz binária original com as 160.000 coordenadas. Arbitramos um valor de  $\alpha = 0,01$  para fins de corte.

## 6 RESULTADOS E DISCUSSÃO

Esta sessão será dividida em filogenia global, filogenia comparativa e desempenho de acordo com as etapas do trabalho. Discutiremos em cada subtópico os aspectos relevantes de cada fase do trabalho. No item 3.1, provamos que se pode gerar grandes filogenias sem exigir muito recurso computacional utilizando o algoritmo SVect. Em 3.2, mostramos e discutimos a comparação feita entre alinhamento e vetorização de sequências. Por fim, no item 3.3, apresentaremos a grande árvore filogenética dos plastídios.

### 6.1. Desempenho

Para demonstrar a vantagem da nossa ferramenta no que diz respeito a tempo de execução, fizemos testes em duas máquinas distintas, cujas configurações encontram-se listadas na tabela I.

Conforme observaremos na próxima seção, os resultados das metodologias contrastadas foram muito similares entre si. Fato que nos levou repetir o teste, aferindo o tempo total de processamento, tanto para o algoritmo SVect, quanto para a abordagem utilizando o

alinhamento múltiplo de sequências. A contagem de tempo computacional contemplou desde o tratamento da sequência bruta até o instante em que o dendrograma é montado, utilizando um computador pessoal.

<b>Tipo de análise</b>	<b>Ordenação (S)*</b>	<b>Espaçadores (s)**</b>	<b>Processamento (s)</b>	<b>Montagem (s)</b>	<b>Total (s)</b>
<b>SVect</b>	-	0,2	<b>0,2</b>	<b>0,4</b>	<b>0,8</b>
<b>Alinhamento</b>	643,59	-	18,2	9,2	670,99

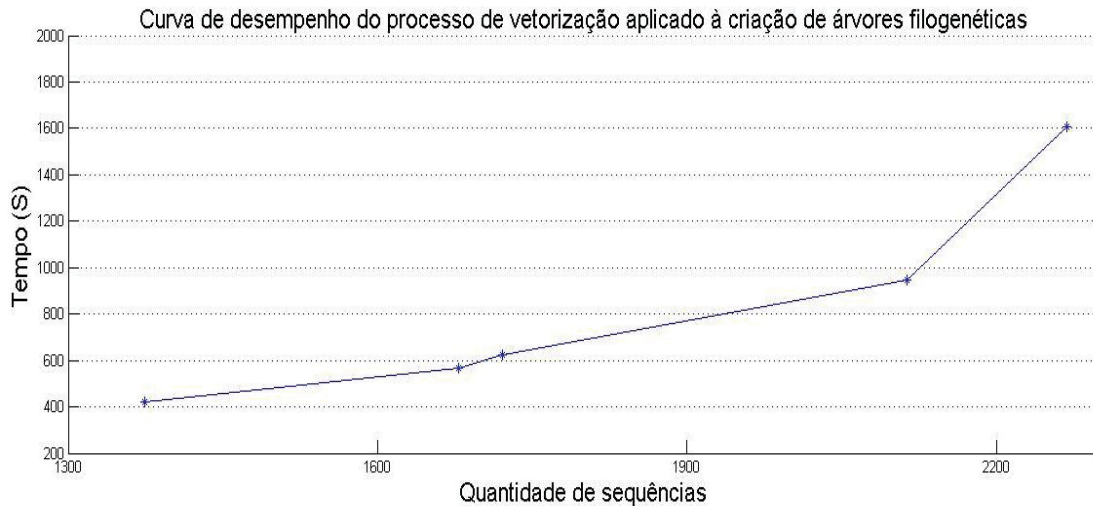
Tabela 2: Tempo de processamento das abordagens de obtenção de árvores filogenéticas comparadas, rodando em um computador pessoal.

Fonte: O Autor. (\*) Na fase de ordenação, as sequências foram concatenadas ordenadamente em todos os organismos, de forma a simular a sintenia genética e favorecer o alinhamento. O SVect dispensa essa ordenação, pois a sintenia não altera os seus resultados. (\*\*) Tempo de pré-processamento para a inserção dos 5 asteriscos como espaçadores. No alinhamento, os espaçadores não são necessários.

Ao considerarmos o tempo de pré-processamento, os números da análise contidos na Tab. 2 nos mostraram que o SVect é cerca de 838 vezes mais rápido quando usamos os dados não projetados. Visto que a ordem de disposição das sequências é essencial para o alinhamento múltiplo de qualidade, tornou-se necessário um subprograma (*script*) em Python3 para ordenar e concatenar todas as sequências no conjunto dos organismos usados como controle. Esta etapa é custosa computacionalmente falando, e como os separadores Inter sequenciais utilizados no SVect permitem o algoritmo gerar os pontos da matriz independentemente em cada proteína, podemos obter a mesma matriz independente da ordem em que as sequências estão dispostas ao longo do proteoma concatenado. Dessa forma, a vetorização de sequências dispensa a etapa de ordenação, o que justifica a diferença gritante de tempo computacional de processamento.

E ainda que desconsiderássemos a etapa de ordenação dos dados pelos algoritmos em questão, o SVect se mantém cerca de 45 vezes mais ágil em relação ao alinhamento múltiplo quando executado nas configurações acima mencionadas.

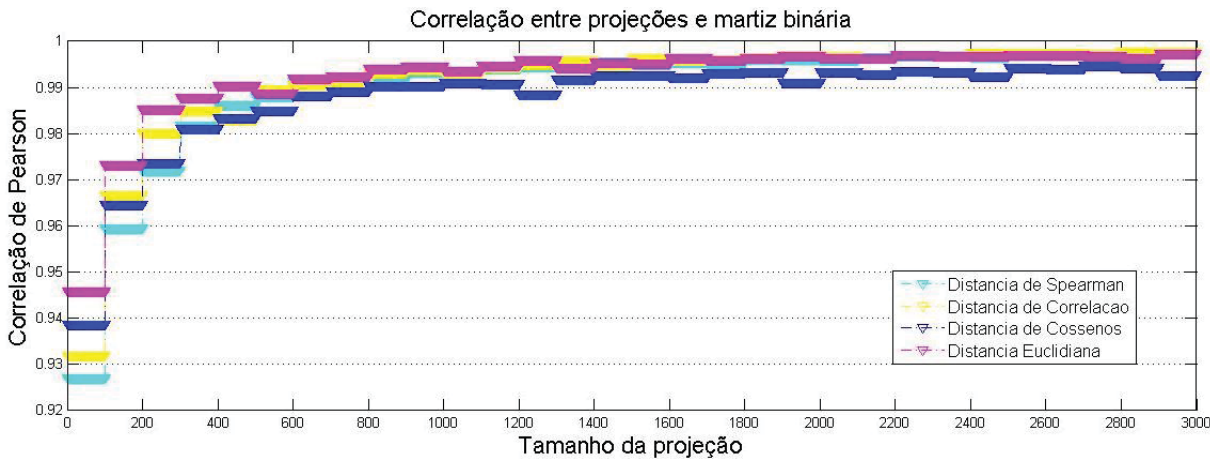
Fizemos também testes para descobrir a curva do tempo de processamento em relação ao número de sequências analisadas sem projeção em um hiper-espço. O resultado está representado na Figura 11, e nos mostra que o crescimento é relativamente linear, mas houve um aumento anormal no tempo de processamento com aproximadamente 2200 sequências, quando houve estouro de memória RAM e o computador passou a utilizar memória do disco rígido para concluir a análise.



**Figura 6:** Gráfico de relação tempo/quantidade de seqüências para obter a árvore filogenética em notebook usando como base o vetor binário não projetado. Fonte: O autor.

Apesar de nos fornecer os melhores resultados, a matriz binária contém 160 mil posições numéricas para cada proteoma/OTU processado. Esta matriz pode tomar muita memória RAM quando muitos organismos forem processados. Logo, a projeção em um hiper-espaço utilizando o teorema Johnsons-Lindenstrauss (JOHNSON; LINDENSTRAUSS, 1984) reduz o tamanho dos dados em memória, aumentando também a velocidade de processamento das informações.

Utilizando 4 algoritmos de cálculo de distância distintos: Spearman, Correlação (Pearson), Cossenos e Euclidiana, projetamos os dados processados pelo SVect em 30 vetores de tamanhos entre 100 e 3000, incrementando de 100 em 100, para inferir a melhor abordagem a ser trabalhada para as projeções em um hiper-espaço usando dados proteômicos de plastídios. Na Figura 12, o gráfico mostra a correlação de Pearson entre as projeções em um hiper-espaço e a matriz binária.



**Figura 7:** Gráfico comparativo dos algoritmos para cálculo de distância em função da correlação entre o tamanho da projeção em um hiper-espaço e a matriz binária. Fonte: O Autor.



A partir da Figura 12 podemos inferir que o algoritmo de distância euclidiana é superior na projeção em um hiper-espço dos dados trabalhados, para a maior parte dos vetores projetados. Isso justifica a escolha desse algoritmo para todas as projeções em um hiper-espço usadas nesse trabalho. Outro detalhe observado é o fato de que todos os vetores projetados a partir do tamanho 600 possuem  $\alpha < 0,01$ . Dessa forma, para sermos parcimoniosos entre a qualidade e o desempenho, usamos e recomendamos a projeção em um hiper-espço de tamanho 600 para trabalhar com dados de proteomas plastidiais.

Tabela 3: Tempo de geração de dendrogramas utilizando diferentes projeções em um hiper-espço em dois computadores distintos.

Projeção em um hiper-espço	Dispositivo	Tempo(S)
<b>Ausente</b>	PC	706,3
	Servidor	393,8
<b>1200</b>	PC	236,6
	Servidor	92,5
<b>1000</b>	PC	217,7
	Servidor	81,8
<b>800</b>	PC	162,6
	Servidor	68,1
<b>600</b>	PC	141,5
	Servidor	59,8
<b>400</b>	PC	125,1
	Servidor	52,8

Fonte: O Autor.

Também contamos o tempo de processamento para construir uma árvore filogenética de 2113 proteomas de plastídios usando o algoritmo SVect a partir do momento em que o arquivo em formato \*.faa é lido. Aferimos o tempo computacional nos dois dispositivos já mencionados para a matriz binária e para as projeções em um hiper-espço com tamanho 1200, 1000, 800, 600 e 400. O resultado foi a obtenção dos cladogramas em tempos que variam entre 52 segundos a pouco mais de 11 minutos (Tab. 3), provando que são viáveis análises por vetorização em computadores pessoais, mesmo com um grande volume de dados.

As 3 tentativas de análise global dos proteomas por alinhamento não obtiveram sucesso pela ferramenta Mega7 (KUMAR et al., 2016), que permaneceu aproximadamente 37 horas

processando e exibiu uma mensagem de erro. O alinhamento de todos os proteomas foi realizado no Omega Clustal (SIEVERS; HIGGINS, 2014) após aproximadamente 45 horas de processamento, mas não foi possível gerar a filogenia, devido a diferença de tamanho dos proteomas alinhados e a má-qualidade do alinhamento por este programa.

## 6.2 Análise Filogenética Comparativa

A etapa amostral do trabalho desenvolvida em Python3®, nos forneceu um conjunto controlado e ordenado de proteínas, tornando possível o alinhamento. Consequentemente, uma árvore baseada na distância de sequências foi obtida. Nesta análise, 15 proteínas comuns e com função conhecida no metabolismo (Quadro 1) foram selecionadas para a comparação, totalizando 21 organismos usados como controle.

Proteínas utilizadas na análise	
<i>photosystem II protein I</i>	<i>ribosomal protein L20</i>
<i>photosystem I subunit VIII</i>	<i>ribosomal protein S12</i>
<i>photosystem II protein J</i>	<i>ribosomal protein S18</i>
<i>photosystem II protein N</i>	<i>cytochrome b6</i>
<i>photosystem II protein D1</i>	<i>cytochrome f</i>
<i>photosystem II protein K</i>	<i>maturase K</i>
<i>photosystem II protein T</i>	<i>ATP synthase CF1 alpha subunit</i>
<i>photosystem I assembly protein</i>	

Quadro 1: Proteínas usadas na análise. Foram utilizadas sequências proteicas essenciais encontradas no maior número de organismos distintos. Resíduos de aminoácidos funcionais na maquinaria de DNA e na fotossíntese compõem a lista. Fonte: O autor.

Ao compararmos o dendrograma gerado por alinhamento e distância entre sequências (Figura 8-A) com o dendrograma irmão construído por nossa metodologia livre de alinhamento (Figura 8-B), verificamos a alteração de dois ramos entre as árvores (*Ludisia discolor* e *Megaleranthus saniculifolia*). Todavia, a topologia não foi alterada, visto que a espécie estão enraizadas na mesma posição em ambos os dendrogramas. Trata-se apenas de diferenças gráficas entre as representações.



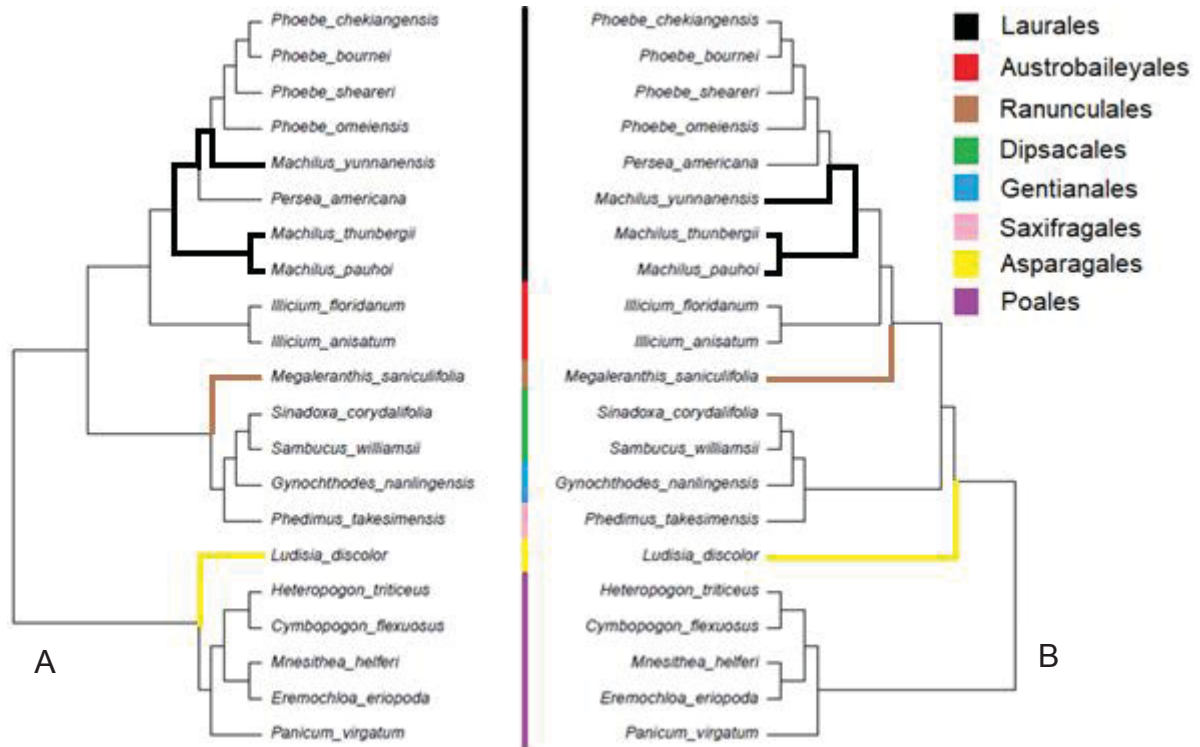


Figura 8: Comparação entre metodologia baseada em alinhamento e livre de alinhamento usando a matriz binária para obtenção de árvores filogenéticas. (A) Topologia construída com base em alinhamento. (B) Topologia gerada com base em vetorização de sequências. Fonte: O autor.

Notamos também que o dendrograma baseado em alinhamento (Fig 8A) dividiu o cluster pertencente ao gênero *Machilus*, inserindo o Abacate (*Persea americana*) entre os integrantes, tornando-o um grupo parafilético. Isso não ocorre no dendrograma baseado em vetorização de sequências (Fig 8B). O gênero *Machilus* coalesce em um único grupo parafilético.

Quando comparamos o dendrograma de baseado em alinhamento, com a árvore filogenética baseada na projeção em um hiper-espaco vetorial de tamanho 800 (Fig. 9), vemos o clado Laurales constituído de mesma topologia. A topologia da árvore muda em *Megaleranthus saniculifolia* (Ranunculales), que é colocado em um grupo monofilético com *Phedimus takesimensis* (Saxifragales). Como ambas as ordens são basais em relação a Dipsacales e Gentianales, não é possível inferir uma topologia como correta. Para tal, seria necessário a inclusão de mais espécies na análise.

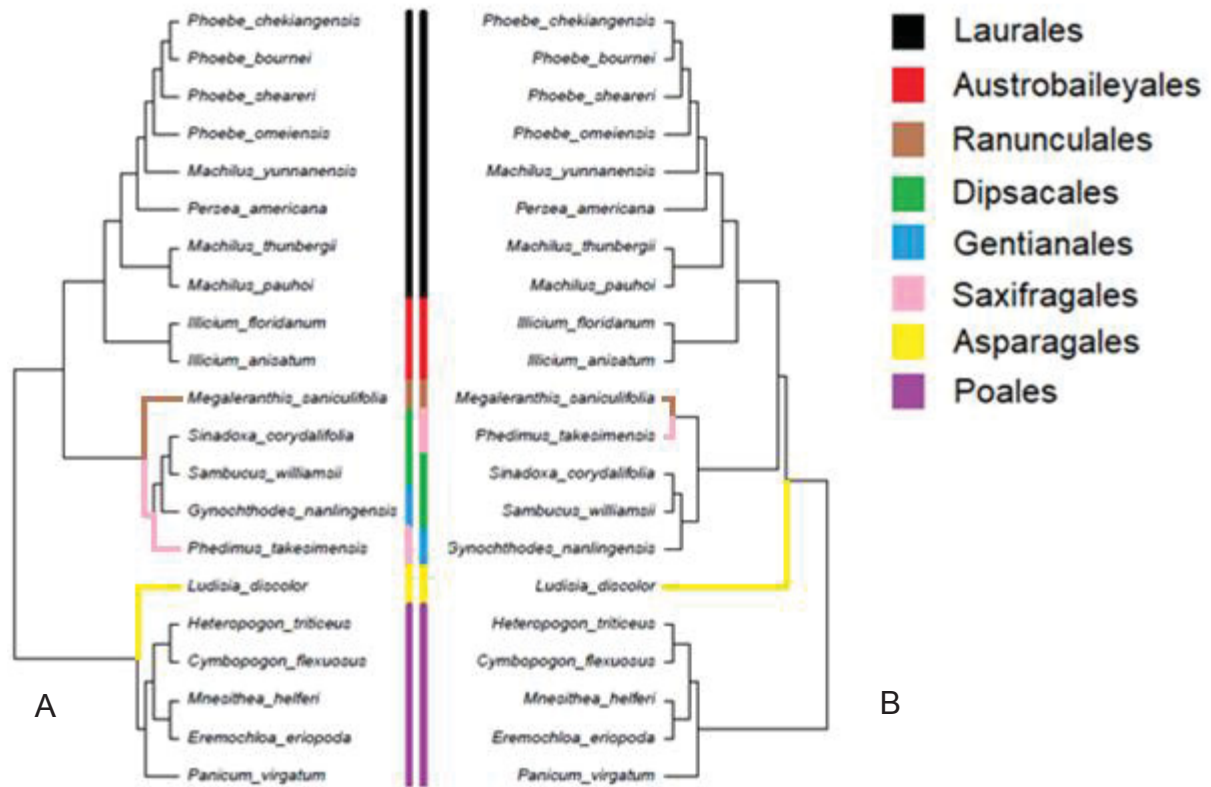


Figura 9: Comparação entre metodologia baseada em alinhamento e livre de alinhamento usando vetor de tamanho 800 para obtenção de árvores filogenéticas. (A) Topologia obtida com base no alinhamento. (B) Topologia construída por vetorização e projeção em um hiper-espço. Fonte: O autor.

O resultado da comparação entre o dendrograma controle obtido por alinhamento e a árvore filogenética gerada pela projeção em um hiper-espço do SVect de tamanho 600 (Fig. 10B) nos dá uma topologia idêntica à do dendrograma gerado pela matriz binária (Fig. 8B). A igualdade entre esses cladogramas é mais uma prova que fundamenta a nossa proposta de trabalho com projeções em um hiper-espço de tamanho 600 no SVect.

Não podemos afirmar aqui que as topologias aqui apresentadas são corretas. Apenas mostramos que nossa abordagem nos fornece resultados muito similares aos do alinhamento múltiplo de sequências, e, em alguns casos, a qualidade da topologia proveniente da vetorização, quando validada na literatura, se mostrou mais consistente, visto que agrupou os plastídios pertencentes ao mesmo gênero (*Machilus*), o que não podemos observar quando utilizamos a metodologia baseada em alinhamento.

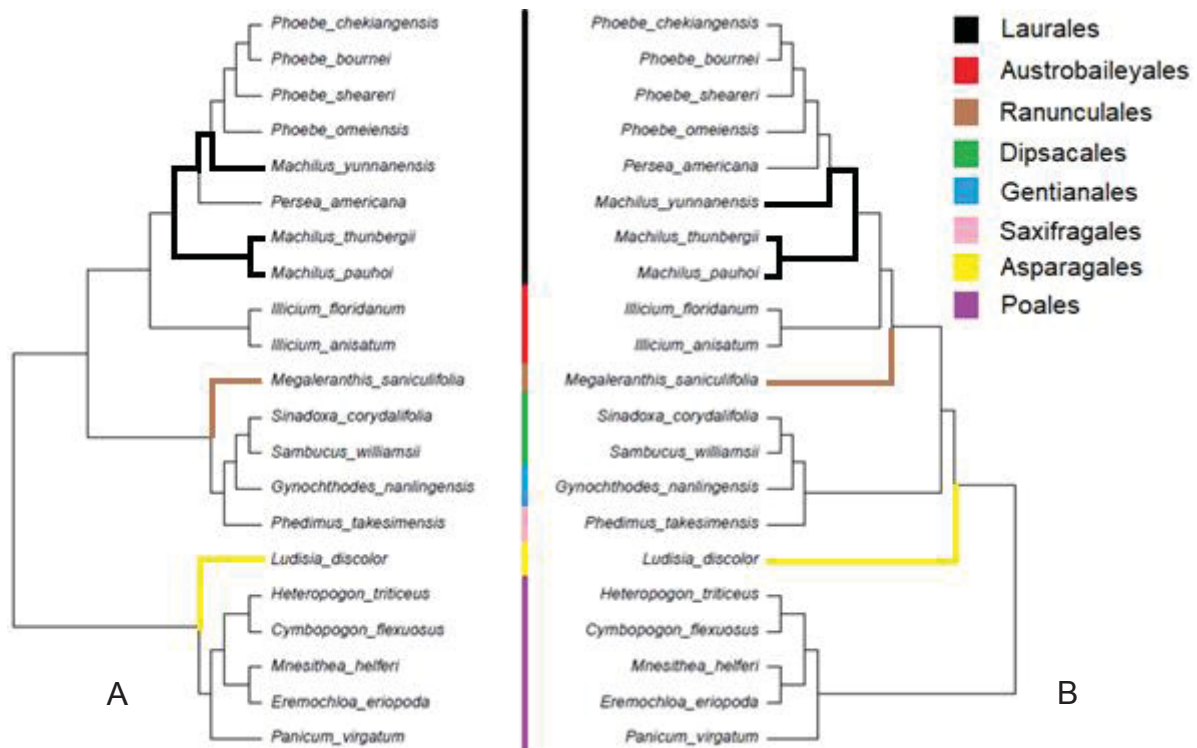


Figura 10: Comparação entre metodologia baseada em alinhamento e livre de alinhamento usando vetor de tamanho 600 para obtenção de árvores filogenéticas (A) Topologia obtida com base no alinhamento. (B) Topologia construída por vetorização e projeção em um hiper-espaço. Fonte: O autor.

### 6.3 Análise Filogenética Global

Baseada no conjunto de dados baixados no RefSeq/Organelle Genome Resource versão 86, disponível no dia 12 de janeiro de 2018, nossa árvore filogenética possui 2493 OTU's, sendo 1962 de espécies portadoras de cloroplastos e 531 portadoras de plastídios indiferenciados e apicoplastos. A análise engloba todos grupos de organismos com plastídios sequenciados e depositados no banco de dados, tais como algas, protozoários, plantas terrestres vasculares e avasculares. Discutiremos partes da árvore em relação à sua coerência, validando os trechos conforme outros trabalhos já publicados. Por razões óbvias, não discutiremos a validação de toda a árvore, todavia, podemos encontrar a filogenia em sua totalidade, de acordo com a topologia sua no Apêndice 1.

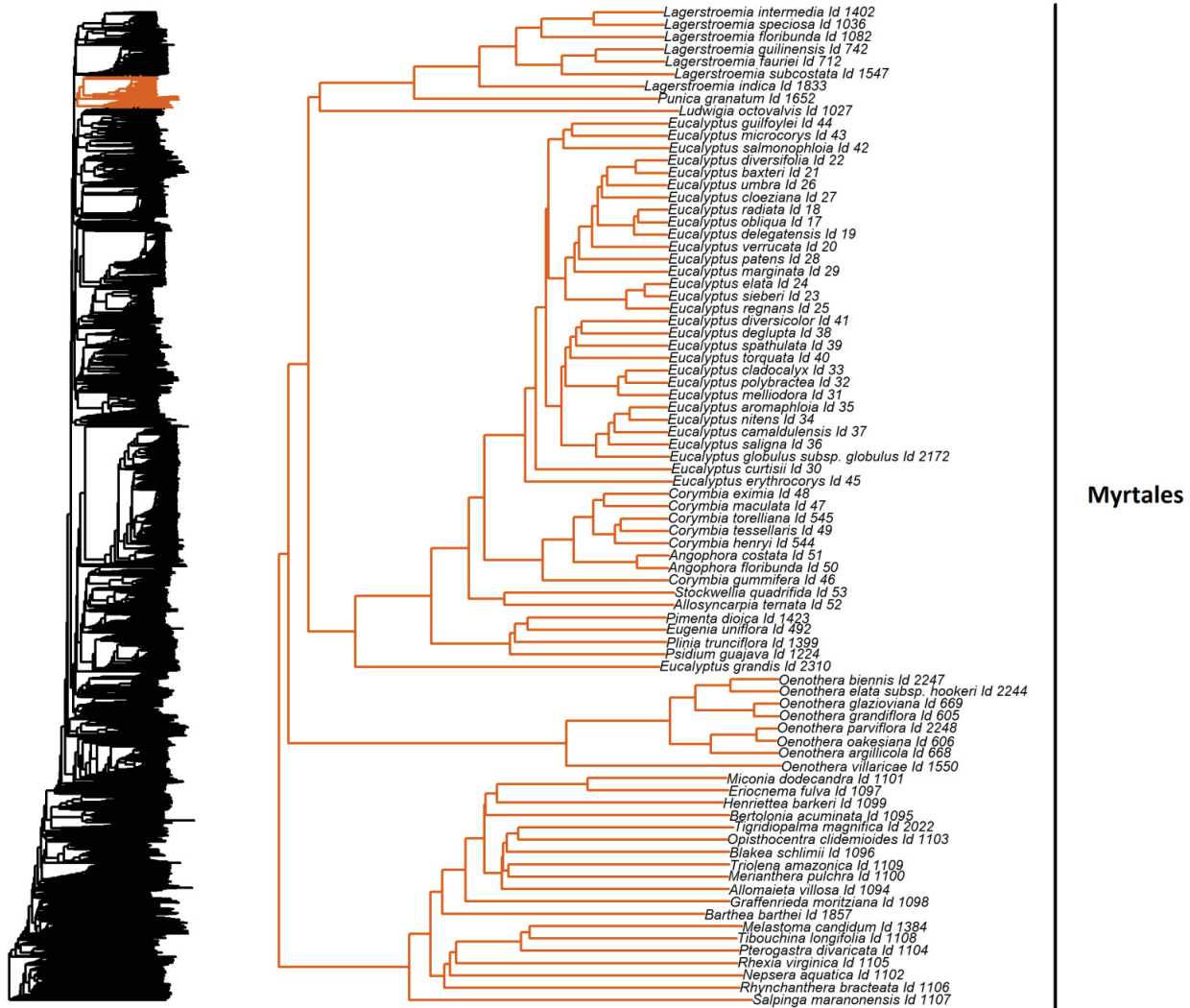


Figura 11: Posição da ordem Myrtales na árvore filogenética vetorial construída pelo método livre de alinhamento SVect. Fonte: O Autor.

A ordem Myrtales é composta por 9 famílias, 380 gêneros, e 13005 espécies, variando entre arbustivas e arbóreas. As folhas são na maioria dos casos opostas e inteiras (sem ramificações, simples), os plastídios encontrados em ductos de condução e armazenam apenas amido, possuem condutores de floema internos e suas flores são geralmente hermafroditas (STEVENS, 2012; BAYLY et al., 2013).

Na Figura 11, os gêneros *Corymbia* e *Angophora* são agrupados juntos. A mesma topologia é confirmada na Figura 12, extraída do estudo de BAYLY et al., (2013), que usa o genoma plastidial completo para inferir a filogenia. Também o gênero *Eucalyptus* em nossa análise encontra-se em similar posição em relação à BAYLY et al., (2013), sendo que os clusters de *E. regnans*/*E. elata*/*E. sieberi* e *E. radiata*/*E. obliqua*/*E. delegatensis*, estão similares em ambos os estudos (Fig. 12). Podemos então inferir que há coerência nos resultados da aplicação do SVect para análise filogenética quando tratamos desta ordem.

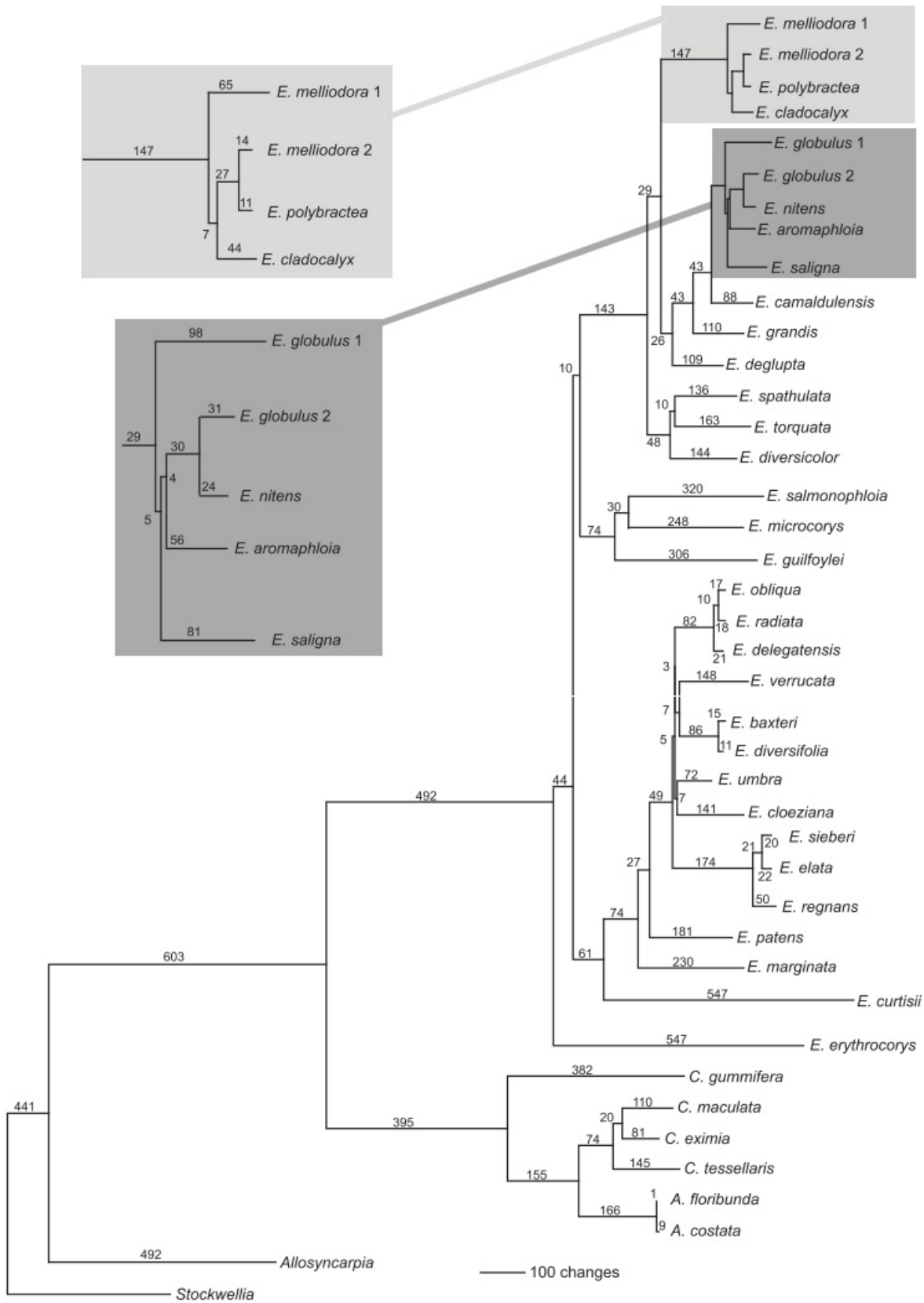


Figura 12: Árvore mais parcimoniosa com base na análise de genomas cloroplastidiais completos de 41 organismos da ordem Myrtales. Fonte: BAYLY et al., (2013).



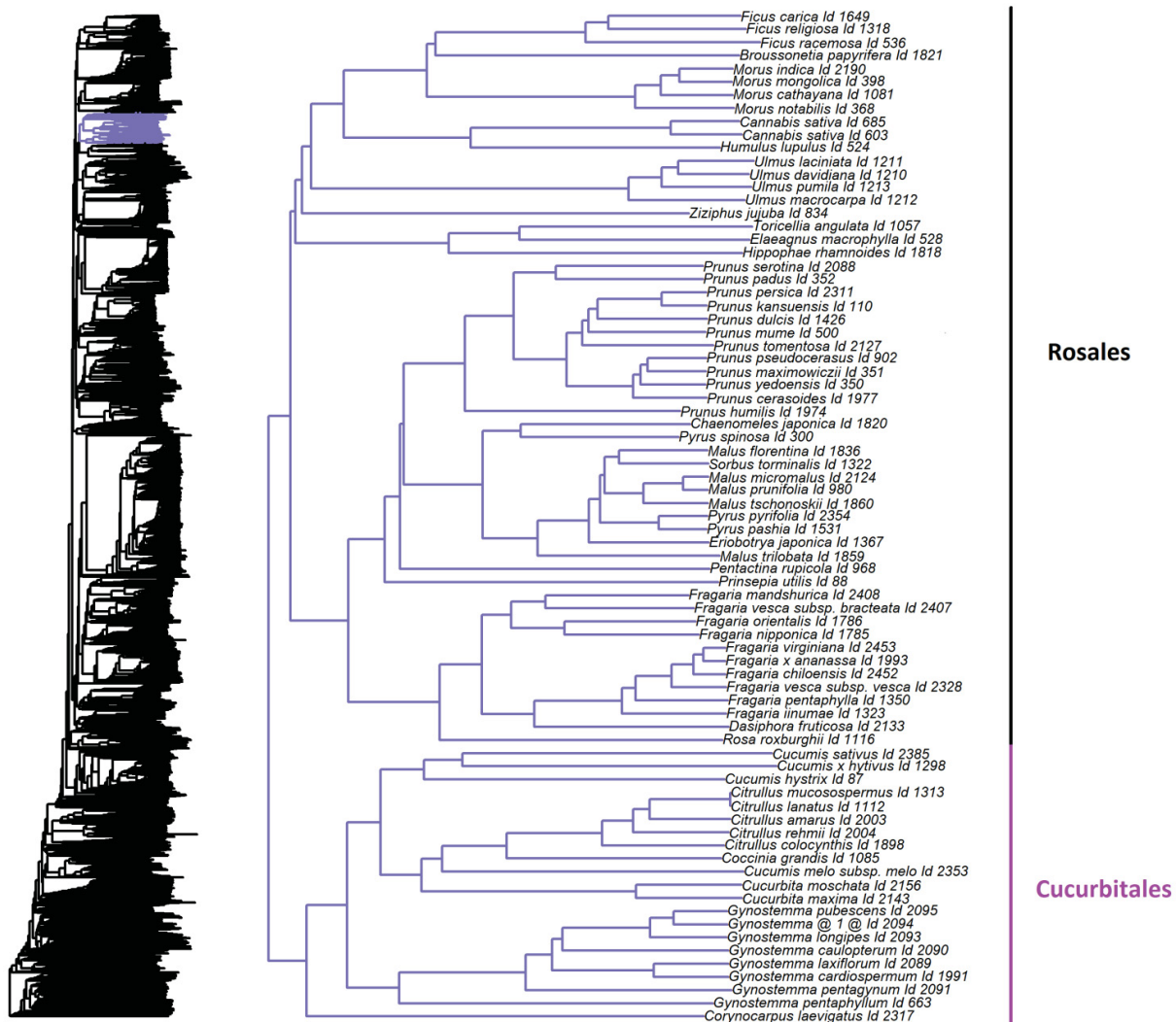


Figura 13: Posição das ordens Rosales e Cucurbitales na árvore filogenética vetorial construída pelo método livre de alinhamento Svect. Trecho da filogenia destacando as ordens em questão. Fonte: O autor.

A ordem das Rosales é composta por 9 famílias, 261 gêneros e 7725 espécies até o presente momento. Suas características sinapomórficas são a ausência de endosperma na semente e a presença de hipanto circundando o ovário (STEVENSON, 2012). É o grupo irmão de Cucurbitales e Fabales, que junto com Fagales, formam o clado de fixação de nitrogênio.

A Figura 13 nos mostra o trecho da filogenia global correspondente à ordem Rosales. Ao compararmos a topologia da Figura 13 com o estudo realizado por LI et al., (2016) vemos as espécies arranjadas de forma similar, seguindo os clados de *Fragaria sp.*, *Prunus sp.* e *Morus sp.*, conforme Figura 14. A topologia também permanece consistente ao comparar com o dendrograma de ZHAO et al., (2017), conforme a Figura 15.

É notório que no trabalho de LI et. al., (2016) o grupo irmão de Rosales é colocado como Cucurbitales (*Corynocarpus* e *Cucumis*), corroborando com nosso estudo.

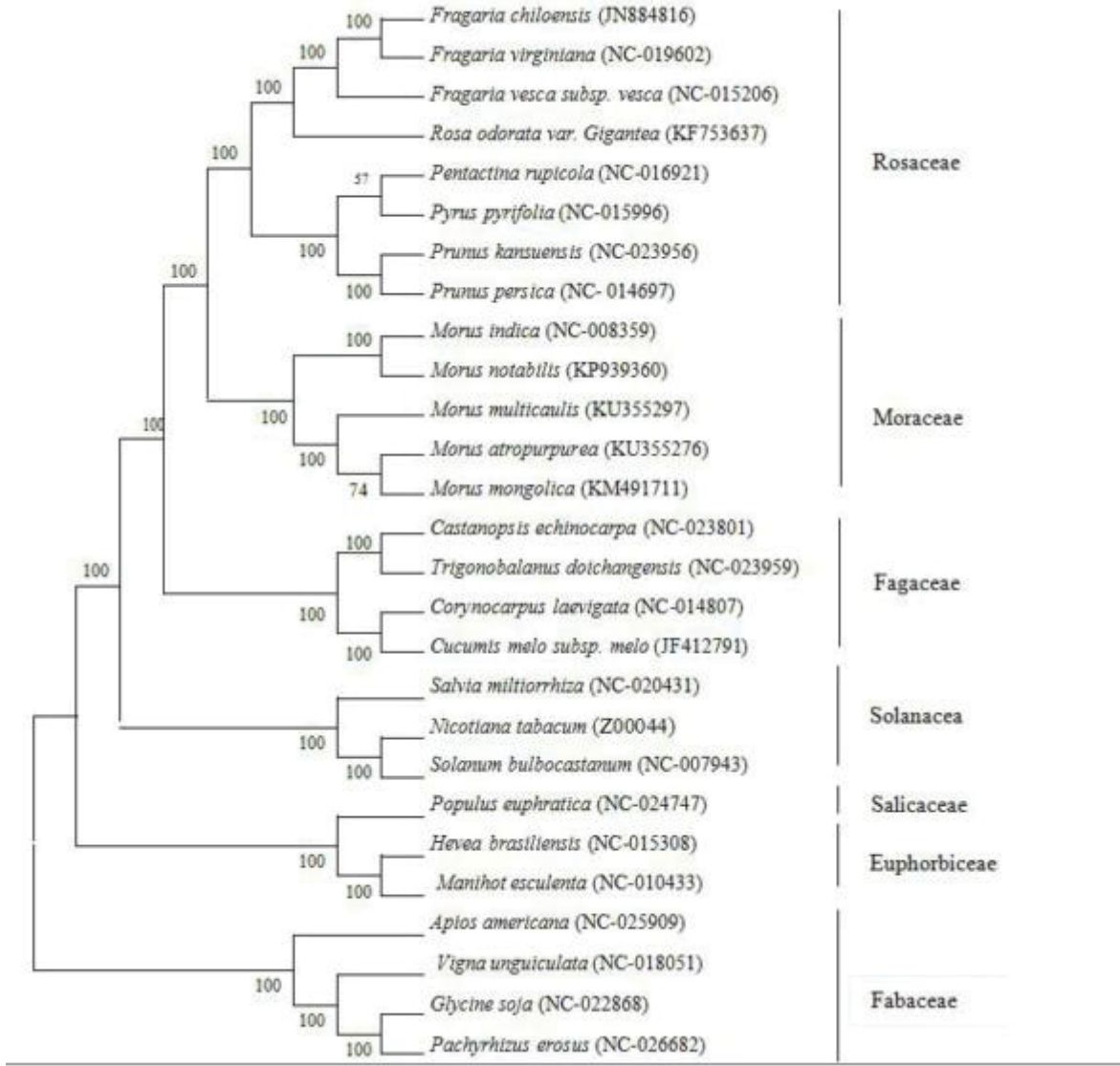


Figura 14: Análise filogenética de espécies de *Amora* baseado em genoma plastidial completo pelo método NJ. Fagaceae, Solanaceae, Salicaceae, Euphorbiceae e fabaceae são incluídos como grupos externos. Fonte:(LI et al., 2016).

Já a ordem Cucurbitales é formada por 7 famílias, 109 gêneros e 2935 espécies e suas principais características são flores unissexuadas, na maioria das vezes pentâmeras e polinizadas por insetos (MATTHEWS; ENDRESS, 2004; STEVENS, 2012).

Em nossa árvore filogenética, a colocação do clado Cucurbitales como grupo irmão das Rosales já foi fundamentado pelo estudo de ZHAO et al., (2017), e a pesquisa de RENNER & SCHAEFER, (2016) fundamenta o clado em nossa árvore colocando *Cucumis*, *Cucurbita*, *Citrullus* e *Coccinia* distante de *Gynostemma*, considerando-o gênero basal, conforme o dendrograma do estudo mostrado na Figura 16.

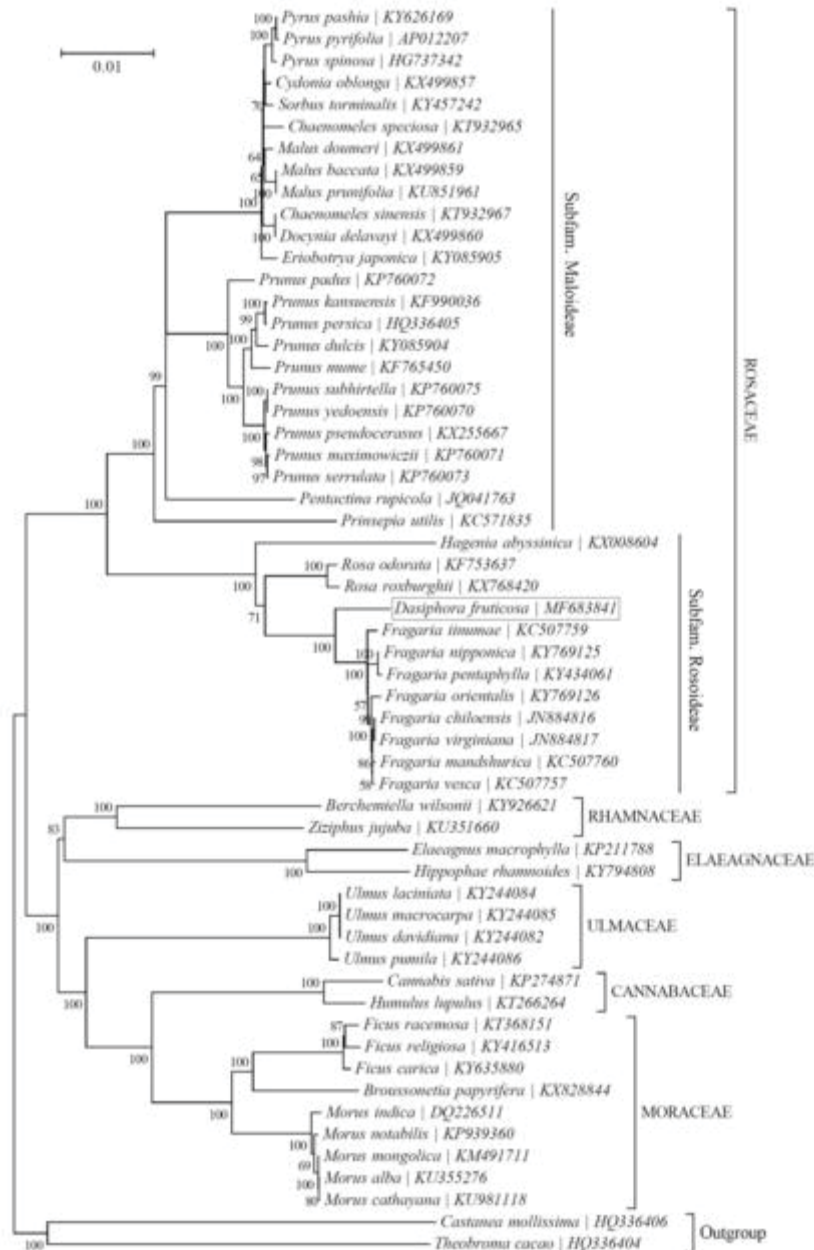


Figura 15: Filogenia de 55 espécies de Rosales gerado a partir de seqüências codificadoras de proteína concatenadas, baseado no algoritmo NJ. os valores de bootstrap são baseados em 1000 replicações. *Castanea mollissima* e *Theobroma cacao* foram incluídos como grupo externo. Fonte: (ZHAO et al., 2017).



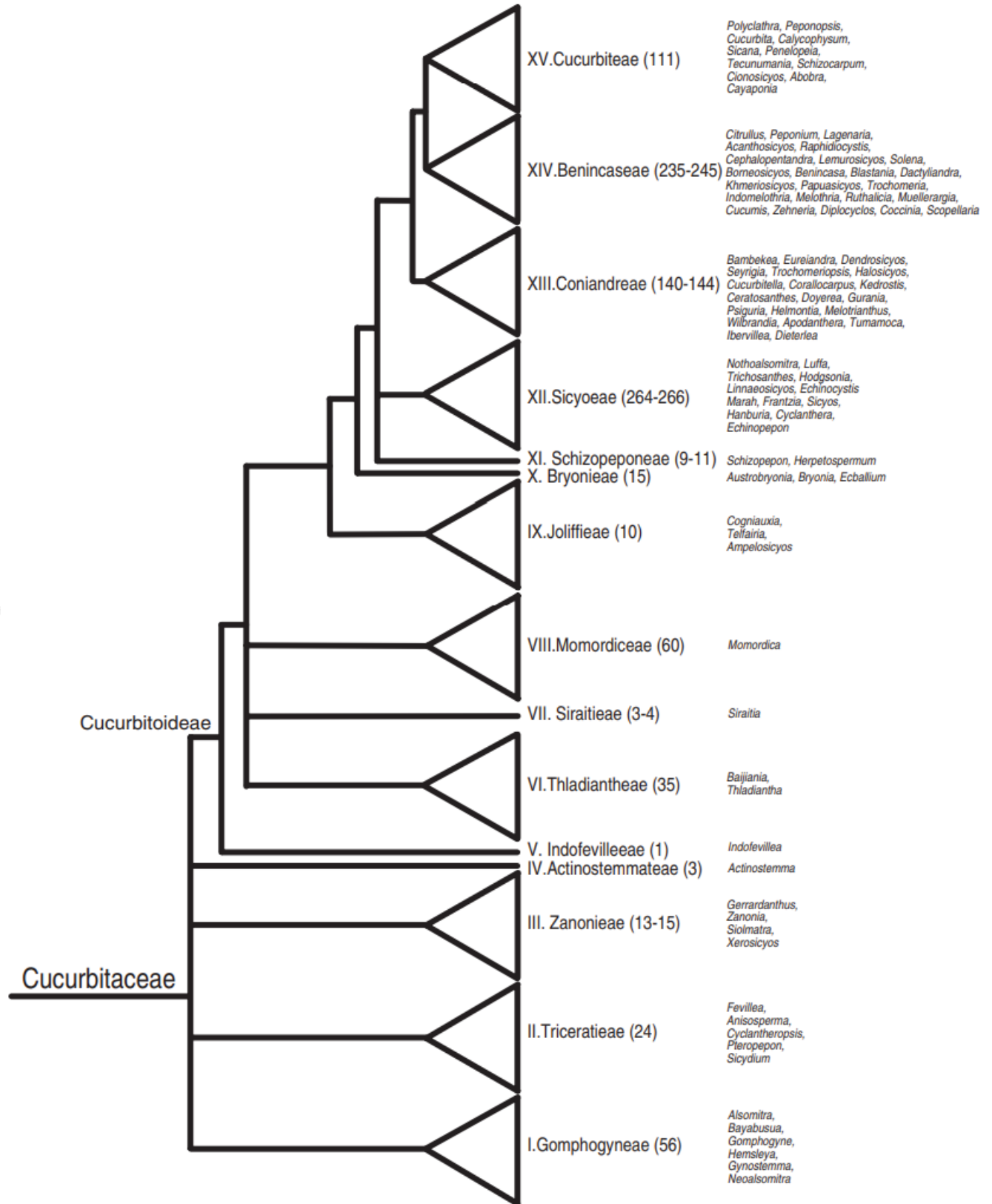


Figura 16: Cladograma baseado em 14 regiões de DNA mitocondrial, nuclear e de cloroplasto, mostrando as relações filogenéticas entre as tribos de Cucurbitaceae. Fonte: (RENNER; SCHAEFER, 2016).

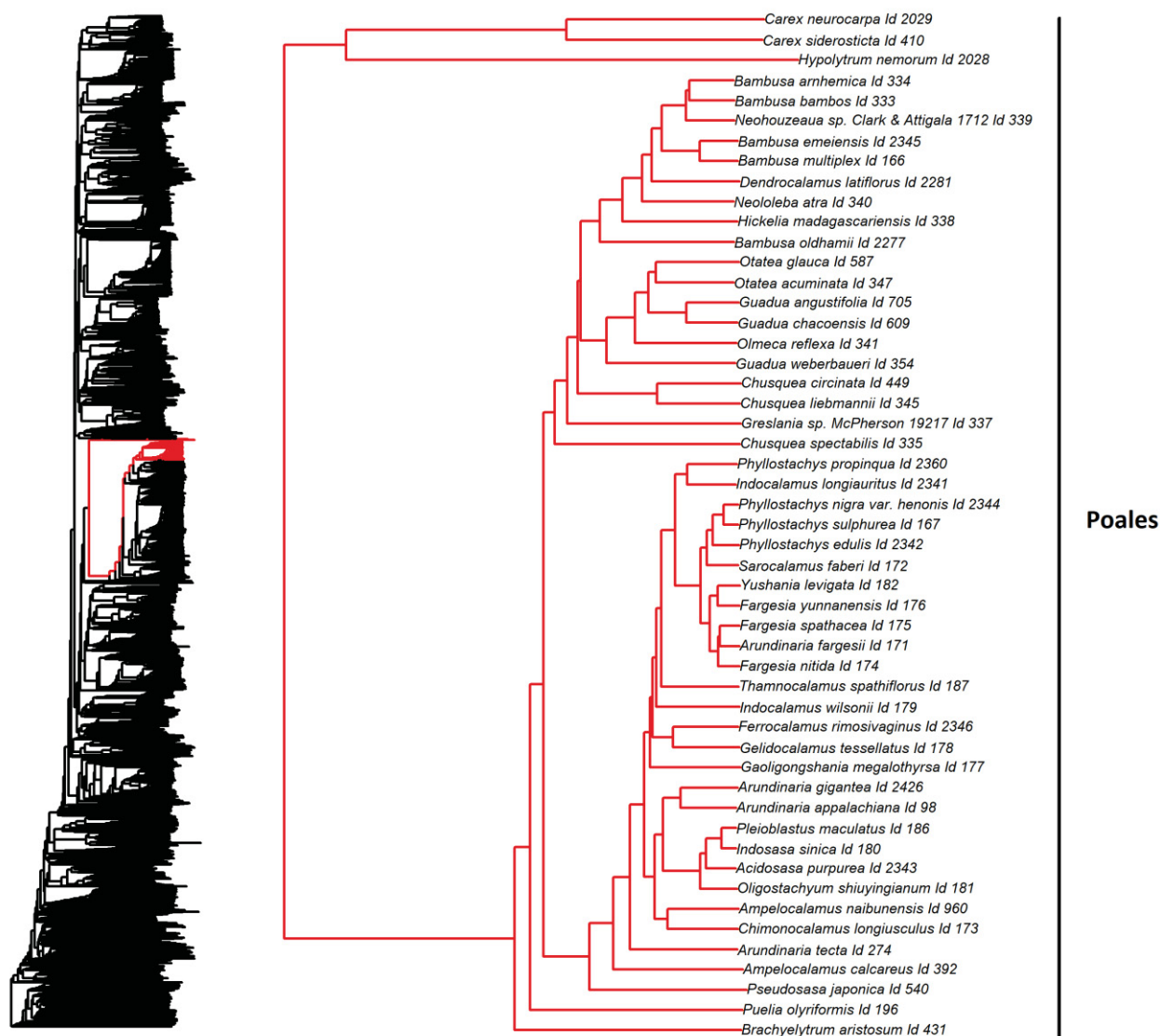


Figura 17: Posição da ordem Poales na árvore filogenética vetorial construída pelo método livre de alinhamento Svect. Trecho da filogenia destacando parte das espécies pertencentes à Poales. Fonte: O autor.

A gigante ordem Poales, contendo 15 famílias, 997 gêneros e 18.875 espécies possui grande importância econômica, pois engloba espécies como o trigo (*Triticum*), o Milho (*Zea*) e o Arroz (*Oryza*). Os espécimes deste táxon possuem pequenas flores revestidas por brácteas dispostas em inflorescências que geralmente são polinizadas por anemocoria, e suas sementes comumente acumulam amido (STEVENSON, 2012).

A Fig. 17 representa uma parte das espécies de Poales utilizadas em nossa análise, e quando sua topologia é comparada com parte da árvore resultante do estudo de SAARELA et al. (2018) representado na Fig. 18, percebemos que os gêneros *Phyllostachis*, *Fargesia*, *Salocalamus*, e *Yushania* estão alocados no grupo irmão de *Arundinaria*, *Indosasa*, *pleioblastus*, e *acidosa*, que por sua vez, estão arranjados como grupo irmão de *Bambusa*, *Greslania*,

*Neoleda Dendrocalamus, Otatea, Olmeca e Guadua*. Assim, validamos a reconstrução filogenética livre de alinhamento também para a ordem Poales.

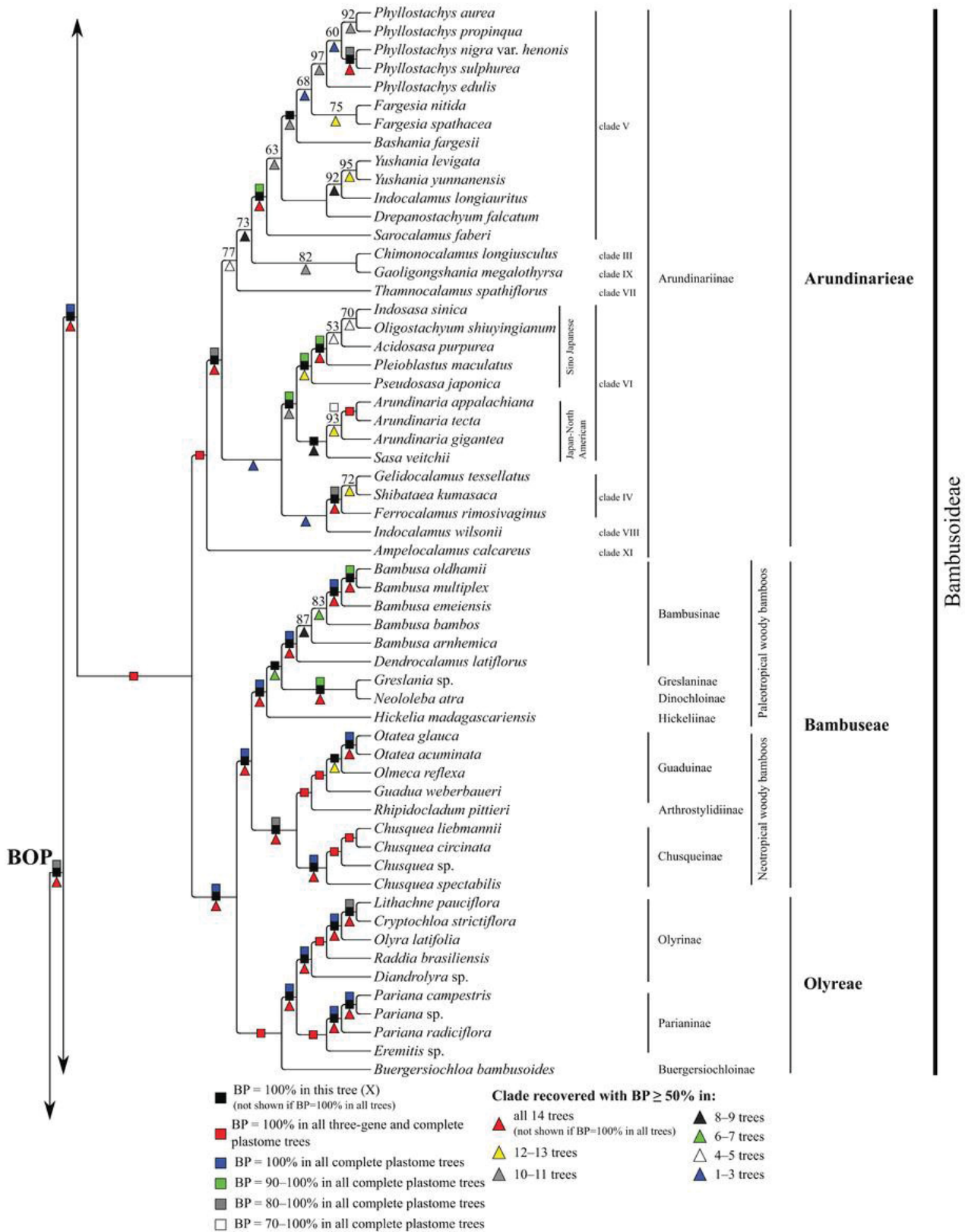


Figura 18: Clado Bambusoideae da árvore baseada em ML e construída por Saarela et. al. (2018) com 250 plastomas completos de espécies da ordem Poales. Fonte: (SAARELA et al., 2018).

A Fig. 19 nos mostra parte do ramo das gimnospermas na grande árvore. As Gimnospermas são um grupo parafilético que inclui as plantas de “sementes nuas”. Pertencem a essa divisão as ordens: Pinales (7 famílias, 68 gêneros, 545 espécies), Gnetales (3 famílias, 3 gêneros, 93 espécies), Gingkoales (1 família, 1 gênero, 1 espécie) e Cycadales (2 famílias, 10 gêneros, 305 espécies), conforme STEVENS (2012).

Ao compararmos a topologia deste grupo na análise global com outros trabalhos, vemos relações muito coerentes. LU et al., (2014) usou os genes nucleares de cópia única LFY e NLY para resolver relações interfamiliares em gimnospermas. A Figura 20 nos permite afirmar que as ramificações e relações entre as ordens diferem sutilmente, mas as topologias a nível de espécie se encontram em consonância. Em Gnetales, por exemplo o clado *Welwitschia/Gnetum/Ephedra*, apesar de diferirem no número de OTU's, possuem um posicionamento filogenético semelhante em relação a árvore global.

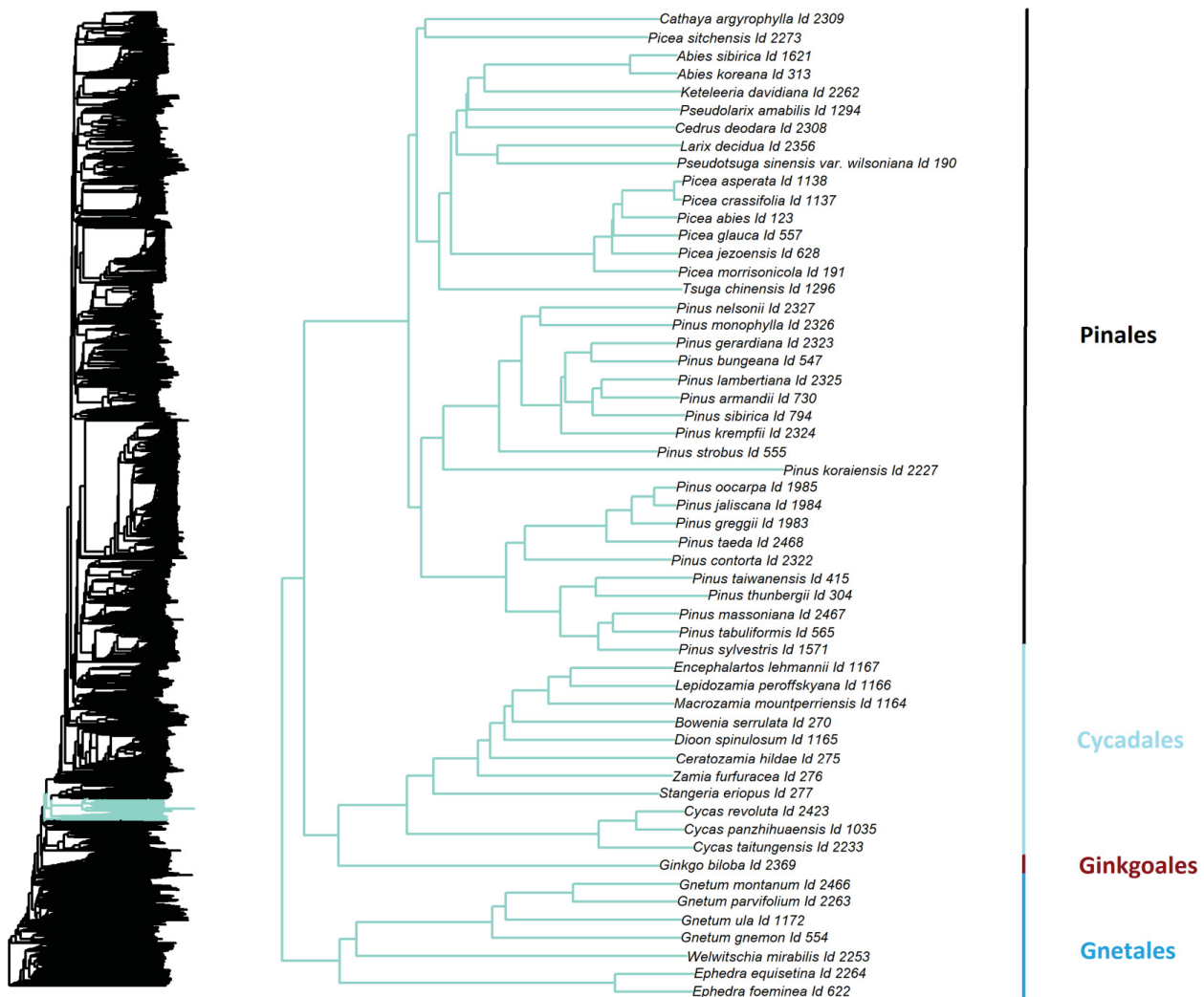


Figura 19: Posição das ordens de Gimnospermas na árvore filogenética vetorial construída pelo método livre de alinhamento SVect. Fonte: O Autor.

Ao compararmos o trecho do dendrograma de LU et al., (2014) equivalente à ordem Cycadales (Fig. 20) com o fragmento correspondente à mesma ordem em nossa filogenia, observamos que ambas as análises colocam o genero *Cycas* como o ramo mais basal do clado, E apesar das diferenças entre o número espécies utilizadas em cada estudo, o nó continua consistente, arranjando organismos estudados pelas duas análises de maneira similar. A mesma situação ocorre também com a ordem Pinales, conferindo confiança ao método e à reconstrução filogenética em questão.

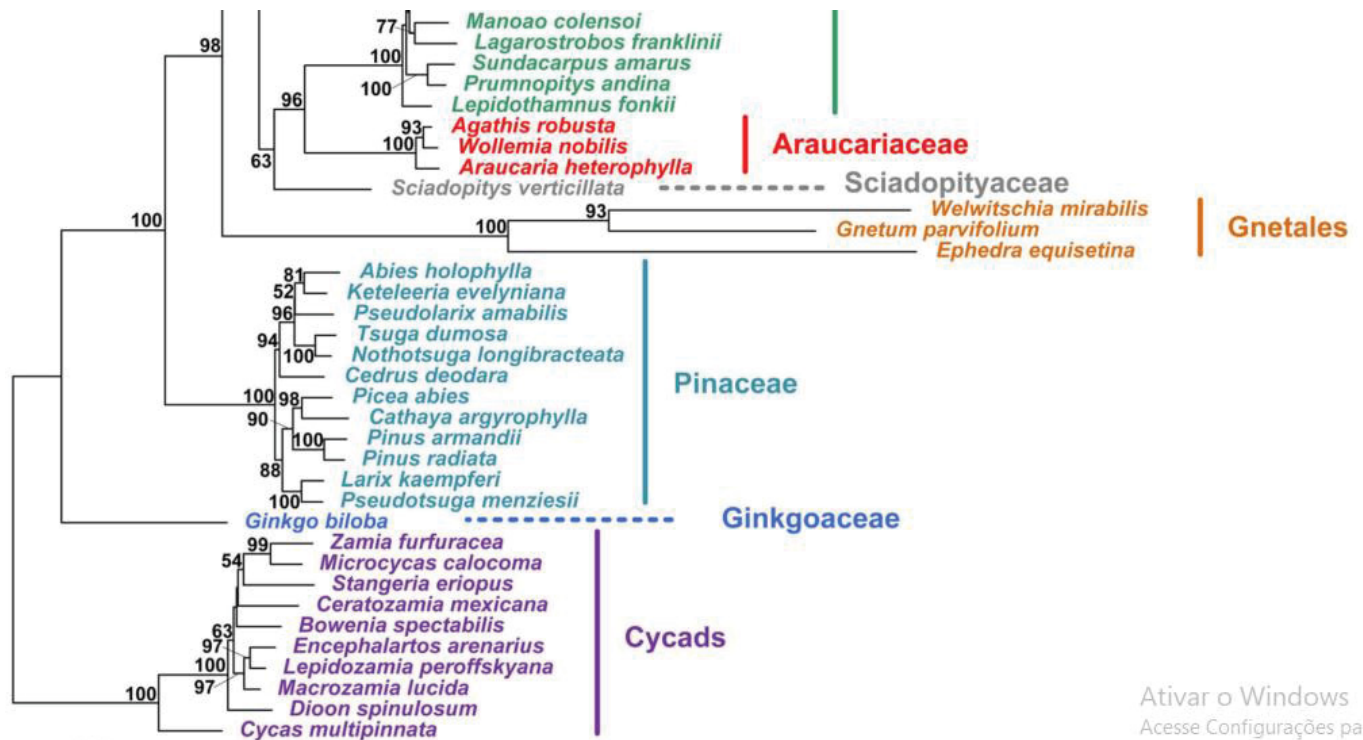


Figura 20: Parte da árvore filogenética das gimnospermas Fonte: Adaptado de LU et al., (2014).

Para mostrar o posicionamento taxonômico de uma nova espécie de *Pinus*, YU et al., (2017) obteve uma análise filogenética Fig (21) que acorda com as posições dos organismos da árvore criada nesse estudo. As espécies *P. taiwanensis*, *P. thunbergii*, *P. tabulaeformis*, *P. massoniana*, *P. contorta* e *P. taeda* se encontram no táxon irmão de *P. koraiensis* e *P. lambertiana*. Assim, podemos afirmar que o arranjo da árvore global de cloroplastos e plastídios para gimnospermas também está em harmonia com outros estudos atuais.



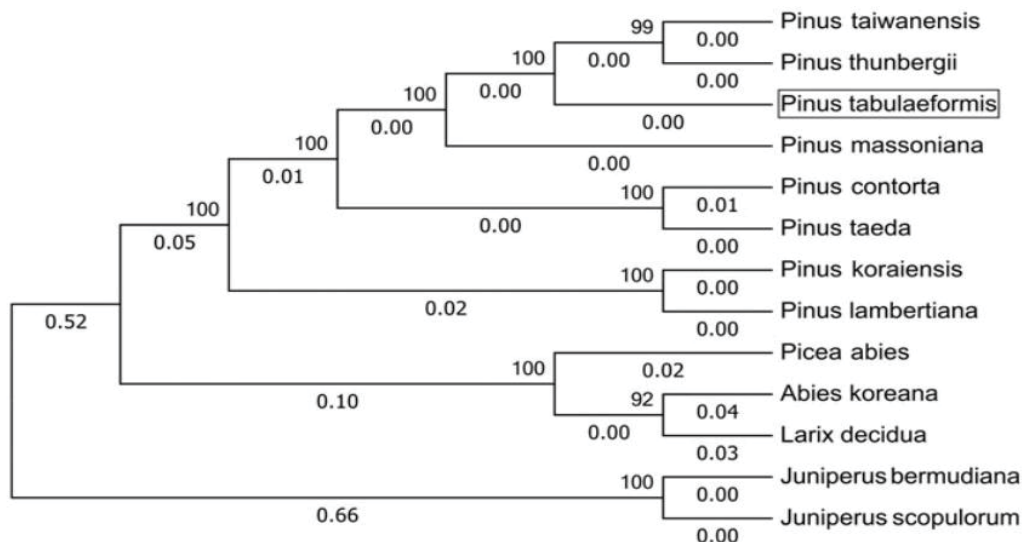


Figura 21: Topologia do cladograma de gimnospermas, com destaque para o gênero *Pinus*. Fonte: (YU et al., 2017).

O fragmento da árvore filogenética global para parte das algas vermelhas é representada pela Figura 22, e encontra respaldo nos científico nos recentes estudos de YANG et al., (2015) e PAIANO et al., (2018). Na Figura 23, YANG et al., (2015) coloca as espécies *Grateloupia taiwanensis* e *Sebdenia flabellata* no mesmo ramo, da mesma forma que o fixa como grupo irmão das Gracilariales, fenômenos também observados na árvore global. As Figuras 22 e 23 ainda acordam em relação à posição do clado Florideophyceae como grupo irmão de Bangiophyceae. Outro fato que confere validade a árvore global é a colocação da classe Cyanidiophyceae como táxon basal da análise, visto que se repete na reconstrução por vetorização. Da mesma forma a Figura 24 retrata uma topologia similar, apesar da diferença no número de espécies.

Diversos estudos também validam a posição das ordens na reconstrução filogenética proposta aqui (MYBURG et al., 2014; RUHFEL et al., 2014; QIAN; ZHANG, 2016; THE ANGIOSPERM PHYLOGENY GROUP VI, 2016; LIU et al., 2017), o que nos permite inferir que o modelo utilizado para as reconstruções filogenéticas são, além de computacionalmente rápidos, eficientes para análises filogenéticas em larga escala baseados em dados proteômicos de plastídios.

Alguns fragmentos da árvore global ficaram indefiníveis. A provável causa dos erros na configuração dos ramos é a má qualidade das sequências. Se os dados de entrada forem mal anotados ou mal sequenciados, toda a análise posterior será prejudicada (PHILIPPE et al., 2017). Além disso, o número de indivíduos de cada táxon presentes na análise tem influência direta sobre a análise filogenética. É necessário uma boa representatividade para uma análise de qualidade.

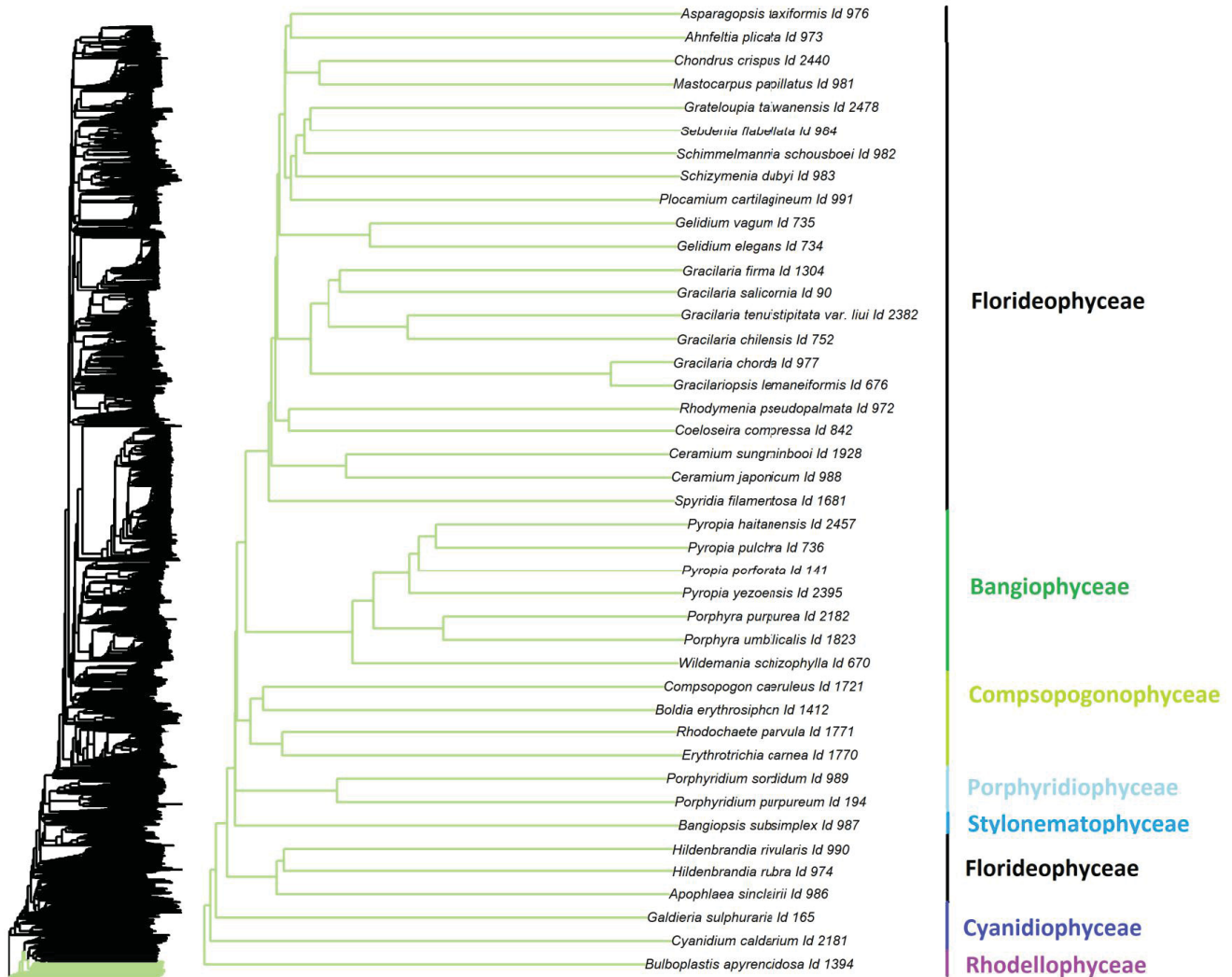


Figura 22: Posição das classes de algas vermelhas na árvore filogenética vetorial construída pelo método livre de alinhamento SVect. Fonte: O Autor.

A validação dos fragmentos da árvore filogenética global foi realizada em sua totalidade por meio da legenda interna de cada figura. As imagens contendo os trechos da árvore foram legendadas por classificação taxonômica de ordem, no caso de plantas vasculares e briófitas, e por classe, no caso das algas. As classificações foram obtidas por meio do *Angiosperm Phylogeny Website* (STEVENS, 2012) e do *Algaebase* (GUIRY; GUIRY, 2018).

Dessa forma, é possível identificar se determinado táxon avaliado tem um certo nível de coerência, pela ordem ou classe do grupo, em comparação às classificações dos organismos subjacentes.

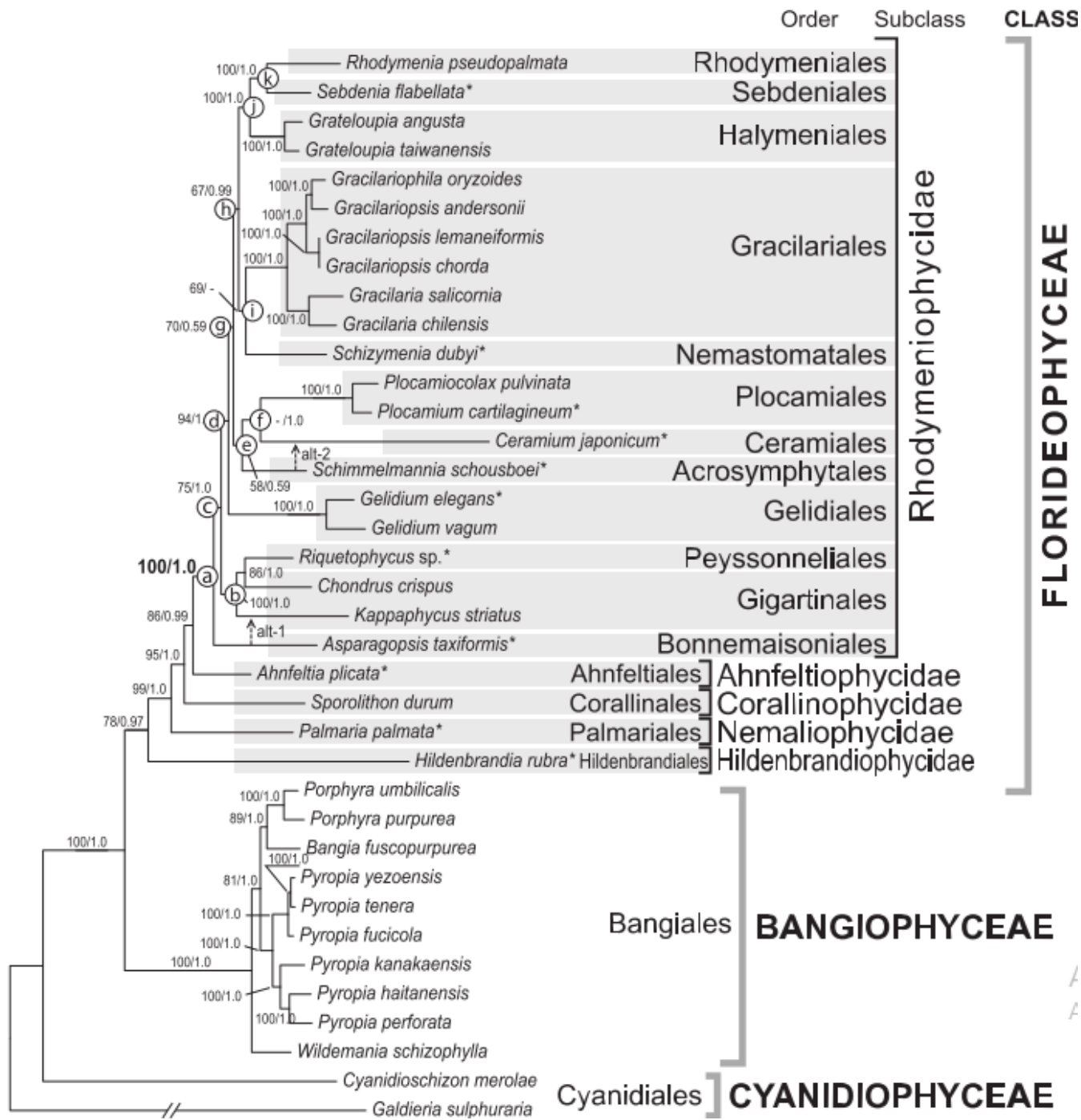
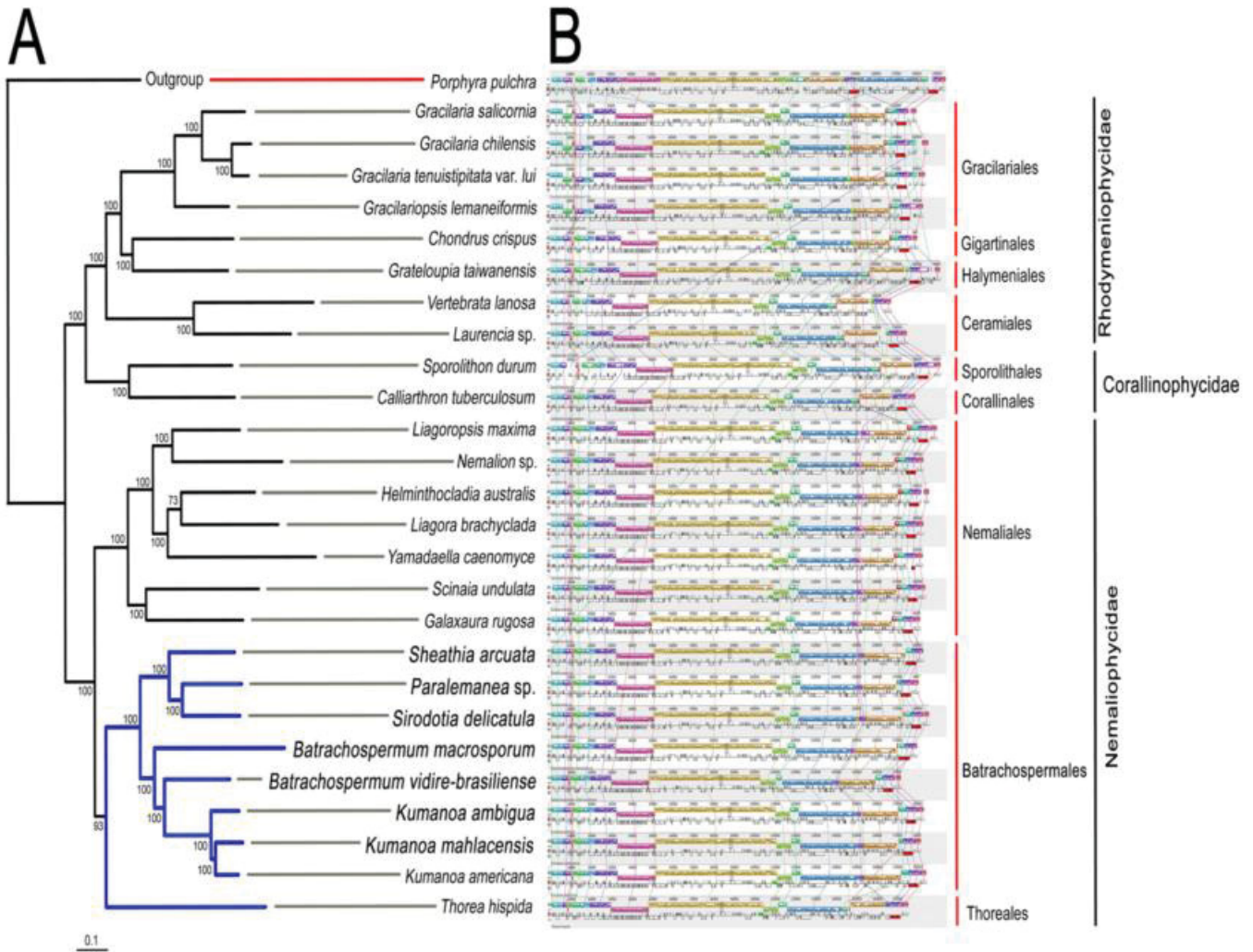


Figura 23: Reconstrução filogenética de algas vermelhas baseada em 24 genes mitocondriais concatenados. Fonte: (YANG et al., 2015)





**Figura 24: Relações filogenéticas e rearranjo plastidial de Florideophyceae, baseados em 177 genes concatenados.**  
**Fonte:** (PAIANO et al., 2018).

## 7. CONCLUSÃO

Por meio desse trabalho, podemos concluir que:

- A filogenômica enfrenta grandes desafios no que diz respeito à utilização do volume de dados disponíveis;
- Tanto a abordagem de super matiz quanto a abordagem de super árvore possuem limitações quando manipulam um grande conjunto de dados, pois utilizam alinhamentos múltiplos;
- O alinhamento múltiplo não atende as demandas atuais de desempenho computacional, e torna-se inviável quando trabalhamos com um conjunto grande de sequências.
- A vetorização de sequências, somado à redução de dimensionalidade mostrou-se viável, rápida e eficiente quando aplicada a análises filogenéticas em larga escala;
- Utilizar a vetorização somada à redução de dimensão em alguns casos pode superar o alinhamento de sequências quando se trata de qualidade de resultados;

- A metodologia proposta se mostrou 45 vezes mais rápida em relação ao alinhamento múltiplo de sequências;
- O SVect torna prescindível a concatenação ordenada para simular uma sintenia genética, o que exprime a vantagem do método em relação ao alinhamento múltiplo de sequências, visto na diferença de tempo de processamento das sequências..
- A qualidade da análise está condicionada a qualidade dos dados de sequência utilizados;
- A distância euclidiana e o tamanho do vetor em 600 pontos numéricos provaram ser o equilíbrio entre qualidade e desempenho para análises envolvendo redução de dimensionalidade em dados de regiões codificantes em genomas plastidiais;
- A árvore filogenética global de plastídios neste trabalho proposta é coerente, e foi validada com os dados já publicados.

## 8. PERSPECTIVAS FUTURAS

O próximo passo é disponibilizar uma ferramenta web para visualização da filogenia global de plastídios proposta nesse trabalho, com a funcionalidade de inserção de novos organismos na árvore. A ferramenta funcionará como classificador, devolvendo ao pesquisador a posição taxonômica dos espécimes cujas CDS's plastidiais sejam submetidas na ferramenta.

Outra necessidade é realizar um estudo detalhado para descobrir o efeito das HTG's (transferências horizontais de genes) nessas árvores resultantes da vetorização e eventual projeção em um hiper-espço.

## REFERÊNCIAS

- ANDERSON, J. M. The molecular organization of chloroplast thylakoids. **BBA Reviews On Bioenergetics**, v. 416, n. 2, p. 191–235, 1975.
- BAYLY, M. J.; RIGAULT, P.; SPOKEVICIUS, A.; et al. Chloroplast genome analysis of Australian eucalypts - *Eucalyptus*, *Corymbia*, *Angophora*, *Allosyncarpia* and *Stockwellia* (Myrtaceae). **Molecular Phylogenetics and Evolution**, v. 69, n. 3, p. 704–716, 2013. Elsevier Inc. Disponível em: <<http://dx.doi.org/10.1016/j.ympev.2013.07.006>>. .
- BENSON, D. A.; CAVANAUGH, M.; CLARK, K.; et al. GenBank. **Nucleic Acids Research**, v. 45, n. D1, p. D37–D42, 2017.
- BLANKENSHIP, R. E. **Molecular Mechanisms of Photosynthesis**. Oxford, UK: Blackwell Scientific, 2013.
- BLEIDORN, C. Third generation sequencing: Technology and its potential impact on evolutionary biodiversity research. **Systematics and Biodiversity**, v. 14, n. 1, p. 1–8, 2016.
- CARVALHO, M. C. D. C. G. DE; SILVA, D. C. G. DA. Sequenciamento de DNA de nova geração e suas aplicações na genômica de plantas. **Ciência Rural**, v. 40, n. 3, p. 735–744, 2010.
- CATANHO, M.; MIRANDA, A. B. DE; DEGRAVE, W. Comparando genomas: bancos de dados e ferramentas computacionais para a análise comparativa de genomas procarióticos. **Reciis**, v. 1, n. 2, p. 335–358, 2007. Disponível em: <<http://www.reciis.cict.fiocruz.br/index.php/reciis/article/view/105/121>>. .
- CHASE, M. W.; SOLTIS, D. E.; OLMSTEAD, R. G.; et al. Phylogenetics of Seed Plants: An Analysis of Nucleotide Sequences from the Plastid Gene *rbcL*. **Annals of the Missouri Botanical Garden**, v. 80, n. 3, p. 528, 1993. Disponível em: <<http://www.jstor.org/stable/2399846?origin=crossref>>. .
- CLEGG, M. T.; GAUT, B. S.; LEARN, G. H.; MORTON, B. R. Rates and patterns of chloroplast DNA evolution. **Proceedings of the National Academy of Sciences of the United States of America**, v. 91, n. 15, p. 6795–67801, 1994.
- CLELAND, W. W.; ANDREWS, T. J.; GUTTERIDGE, S.; HARTMAN, F. C.; LORIMER, G. H. Mechanism of Rubisco: The Carbamate as General Base <sup>X</sup>. **Chemical Reviews**, v. 98, n. 2, p. 549–562, 1998. Disponível em: <<http://pubs.acs.org/doi/abs/10.1021/cr970010r>>. .
- CUI, L.; VEERARAGHAVAN, N.; RICHTER, A.; et al. ChloroplastDB: the Chloroplast Genome Database. **Nucleic acids research**, v. 34, n. Database issue, p. D692–6, 2006. Disponível em: <[http://nar.oxfordjournals.org/content/34/suppl\\_1/D692.abstract](http://nar.oxfordjournals.org/content/34/suppl_1/D692.abstract)>. Acesso em: 14/6/2016.
- DANIELL, H.; LIN, C.-S.; YU, M.; CHANG, W.-J. Chloroplast genomes: diversity, evolution, and applications in genetic engineering. **Genome Biology**, v. 17, n. 1, p. 134, 2016. Genome Biology. Disponível em: <<http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-1004-2>>. .
- DAY, W. H. E. Computational complexity of inferring phylogenies from dissimilarity matrices. **Bulletin of Mathematical Biology**, v. 49, n. 4, p. 461–467, 1987.
- DINIZ, F. Biotecnologia : modernas técnicas de engenharia genética unem setores em prol do desenvolvimento nacional. **Biotecnologia Ciência & Desenvolvimento**, , n. 37, p. 2–5, 2007.
- EGEA, I.; BARSAN, C.; BIAN, W.; et al. Chromoplast differentiation: Current status and perspectives. **Plant and Cell Physiology**, v. 51, n. 10, p. 1601–1611, 2010.





representing 19 angiosperm families and one gymnosperm family based on 390 orthologous genes. **Plant Systematics and Evolution**, v. 303, n. 3, p. 413–417, 2017. Springer Vienna.

LOPEZ-JUEZ, E.; PYKE, K. A. Plastids unleashed: Their development and their integration in plant development. **International Journal of Developmental Biology**, v. 49, n. 5–6, p. 557–577, 2005.

LU, Y.; RAN, J. H.; GUO, D. M.; YANG, Z. Y.; WANG, X. Q. Phylogeny and divergence times of gymnosperms inferred from single-copy nuclear genes. **PLoS ONE**, v. 9, n. 9, 2014.

LYUBETSKY, .V A.; SELIVERSTOV, A. V.; ZVERKOV, O. A. Построение разделяющих паралоги семейств гомологичных белков, кодируемых в пластидах цветковых растений. **Mathematical Biology and Bioinformatics**, v. 8, n. 1, p. 225–233, 2013.

MATTHEWS, M. L.; ENDRESS, P. K. Comparative floral structure and systematics in Cucurbitales (Corynocarpaceae, Coriariaceae, Tetramelaceae, Datisceae, Begoniaceae, Cucurbitaceae, Anisophylleaceae). **Botanical Journal of the Linnean Society**, v. 145, n. 2, p. 129–185, 2004. Disponível em:  
<<http://springerlink.metapress.com/openurl.asp?genre=article&id=doi:10.1007/s00606-003-0090-2>>. .

MCCARTHY, C. G. P.; FITZPATRICK, D. A. Phylogenomic Reconstruction of the Oomycete Phylogeny Derived from 37 Genomes. **mSphere**, v. 2, n. 2, p. 1–17, 2017. Disponível em: <[http://msphere.asm.org/content/2/2/e00095-17?utm\\_source=TrendMDmSphere&utm\\_medium=TrendMDmSphere&utm\\_campaign=trendmdalljournals](http://msphere.asm.org/content/2/2/e00095-17?utm_source=TrendMDmSphere&utm_medium=TrendMDmSphere&utm_campaign=trendmdalljournals)>. .

MCMAHON, M. M.; SANDERSON, M. J. Phylogenetic supermatrix analysis of GenBank sequences from 2228 papilionoid legumes. **Systematic Biology**, v. 55, n. 5, p. 818–836, 2006.

MYBURG, A. A.; GRATTAPAGLIA, D.; TUSKAN, G. A.; et al. The genome of *Eucalyptus grandis*. **Nature**, v. 510, n. 7505, p. 356–362, 2014.

NICOLAU, P. B. História Da Classificação Biológica. , 2017. Disponível em:  
<[https://repositorioaberto.uab.pt/bitstream/10400.2/6133/1/texto\\_apoio\\_1\\_Hist\\_classif\\_biologica.pdf](https://repositorioaberto.uab.pt/bitstream/10400.2/6133/1/texto_apoio_1_Hist_classif_biologica.pdf)>. Acesso em: 7/3/2017.

O'BRIEN, E. A.; ZHANG, Y.; WANG, E.; et al. GOBASE: An organelle genome database. **Nucleic Acids Research**, v. 37, n. SUPPL. 1, p. 946–950, 2009.

O'LEARY, N. A.; WRIGHT, M. W.; BRISTER, J. R.; et al. Reference sequence (RefSeq) database at NCBI: Current status, taxonomic expansion, and functional annotation. **Nucleic Acids Research**, v. 44, n. D1, p. D733–D745, 2016.

PAIANO, M. O.; CORTONA, A. DEL; COSTA, J. F.; et al. Organization of plastid genomes in the freshwater red algal order Batrachospermales (Rhodophyta). **Journal of Phycology**, v. 54, n. 1, p. 25–33, 2018.

PATWARDHAN, A.; RAY, S.; ROY, A. Biology Molecular Markers in Phylogenetic Studies-A Review. **J Phylogen Evolution Biol**, v. 2, n. 22, 2014. Disponível em: <<http://dx.doi.org/10.4172/2329-9002.1000131>>. Acesso em: 9/3/2017.

PHILIPPE, H.; VIENNE, D. M. DE; RANWEZ, V.; et al. Pitfalls in supermatrix phylogenomics. **European Journal of Taxonomy**, , n. 283, p. 1–25, 2017. Disponível em: <<http://www.europeanjournaloftaxonomy.eu/index.php/ejt/article/view/407>>. .

PIERRI, C. R. DE. **REPRESENTAÇÕES VETORIAIS DE PROTEOMAS : UM ESTUDO DE CASO REPRESENTAÇÕES VETORIAIS DE PROTEOMAS : UM ESTUDO DE CASO**, 2017. Universidade Federal do Paraná.

- PRUITT, K.; BROWN, G.; TATUSOVA, T.; MAGLOTT, D. Chapter 18 : The Reference Sequence (RefSeq) Database Database. **The NCBI Handbook**, , n. Bethesda(MD): National Center for Biotechnology Information (US), p. 1–22, 2002.
- PRUITT, K. D.; TATUSOVA, T.; MAGLOTT, D. R. NCBI reference sequences (RefSeq): A curated non-redundant sequence database of genomes, transcripts and proteins. **Nucleic Acids Research**, v. 35, n. SUPPL. 1, p. 501–504, 2007.
- PYRON, R.; BURBRINK, F. T.; WIENS, J. J.; et al. A phylogeny and revised classification of Squamata, including 4161 species of lizards and snakes. **BMC Evolutionary Biology**, v. 13, n. 1, p. 93, 2013. Disponível em: <<http://bmcevolbiol.biomedcentral.com/articles/10.1186/1471-2148-13-93>>. .
- QIAN, H.; ZHANG, J. Are phylogenies derived from family-level supertrees robust for studies on macroecological patterns along environmental gradients? **Journal of Systematics and Evolution**, v. 54, n. 1, p. 29–36, 2016.
- RANNALA, B.; YANG, Z. H. Probability distribution of molecular evolutionary trees: A new method of phylogenetic inference. **Journal of Molecular Evolution**, v. 43, n. 3, p. 304–311, 1996.
- REITH, M.; MUNHOITAND, J. Complete Nucleotide Sequence of the Porphyra purpurea Chloroplast Genome. **Plant Molecular Biology Reporter**, v. 13, n. 4, p. 333–335, 1995.
- RENNER, S. S.; SCHAEFER, H. Phylogeny and Evolution of the Cucurbitaceae. **Plant Genetics and Genomics: Crops and Models**, p. 1–26, 2016. Disponível em: <[springer.com/chapter/10.1007/7397\\_2016\\_14](http://springer.com/chapter/10.1007/7397_2016_14)>. .
- ROGALSKI, M.; NASCIMENTO VIEIRA, L. DO; FRAGA, H. P.; GUERRA, M. P. Plastid genomics in horticultural species: importance and applications for plant population genetics, evolution, and biotechnology. **Frontiers in Plant Science**, v. 6, n. July, p. 1–17, 2015. Disponível em: <[http://www.frontiersin.org/Crop\\_Science\\_and\\_Horticulture/10.3389/fpls.2015.00586/abstract](http://www.frontiersin.org/Crop_Science_and_Horticulture/10.3389/fpls.2015.00586/abstract)>. .
- RUHFEL, B. R.; GITZENDANNER, M. A.; SOLTIS, P. S.; et al. From algae to angiosperms—inferring the phylogeny of green plants (Viridiplantae) from 360 plastid genomes. **BMC Evolutionary Biology**, v. 14, n. 1, p. 23, 2014. Disponível em: <<http://bmcevolbiol.biomedcentral.com/articles/10.1186/1471-2148-14-23>>. .
- SAARELA, J. M.; BURKE, S. V.; WYSOCKI, W. P.; et al. A 250 plastome phylogeny of the grass family (Poaceae): topological support under different data partitions. **PeerJ**, v. 6, p. e4299, 2018. Disponível em: <<https://peerj.com/articles/4299>>. .
- SABLOK, G.; MUDUNURI, S. B.; PATNANA, S.; et al. Chloromitosrdb: Open source repository of perfect and imperfect repeats in organelle genomes for evolutionary genomics. **DNA Research**, v. 20, n. 2, p. 127–133, 2013.
- SABLOK, G.; PADMA RAJU, G. V.; MUDUNURI, S. B.; et al. ChloroMitoSSRDB 2.00: More genomes, more repeats, unifying SSRs search patterns and on-the-fly repeat detection. **Database**, v. 2015, p. 1–10, 2015.
- SAITOU, N.; NEI, M. The neighbor-joining method: a new method for reconstructing phylogenetic trees. **Molecular Biology and Evolution**, , n. January, 1987. Disponível em: <<https://academic.oup.com/mbe/article/4/4/406/1029664/The-neighborjoining-method-a-new-method-for>>. .
- SANGER, F.; COULSON, A. R. A rapid method for determining sequences in DNA by primed

synthesis with DNA polymerase. **Journal of Molecular Biology**, v. 94, n. 3, p. 441–448, 1975. Disponível em: <<http://linkinghub.elsevier.com/retrieve/pii/0022283675902132>>. .

SHANKER, A. ChloroSSRdb: a repository of perfect and imperfect chloroplastic simple sequence repeats (cpSSRs) of green plants. **Database : the journal of biological databases and curation**, v. 2014, n. 2, p. 1–5, 2014.

SHAW, J.; LICKEY, E. B.; SCHILLING, E. E.; SMALL, R. L. Comparison of whole chloroplast genome sequences to choose noncoding regions for phylogenetic studies in angiosperms: The Tortoise and the hare III. **American Journal of Botany**, v. 94, n. 3, p. 275–288, 2007.

SHENDURE, J.; BALASUBRAMANIAN, S.; CHURCH, G. M.; et al. DNA sequencing at 40: Past, present and future. **Nature**, v. 550, n. 7676, 2017. Nature Publishing Group. Disponível em: <<http://dx.doi.org/10.1038/nature24286>>. .

SHENDURE, J.; MITRA, R. D.; VARMA, C.; CHURCH, G. M. Advanced sequencing technologies: methods and goals. **Nature reviews. Genetics**, v. 5, n. 5, p. 335–344, 2004. Disponível em: <<http://www.nature.com/nrg/journal/v5/n5/full/nrg1325.html>>. .

SIEVERS, F.; HIGGINS, D. G. Clustal Omega. **Current Protocols in Bioinformatics**, v. 2014, n. December, p. 3.13.1–3.13.16, 2014.

SIM, A. Y. L.; MINARY, P.; LEVITT, M. Modeling nucleic acids. **Current Opinion in Structural Biology**, v. 22, n. 3, p. 273–278, 2012. Elsevier Ltd. Disponível em: <<http://dx.doi.org/10.1016/j.sbi.2012.03.012>>. .

SIMPSON, J. T.; POP, M. The Theory and Practice of Genome Sequence Assembly. **Annu. Rev. Genomics Hum. Genet**, v. 16, p. 153–72, 2015.

SMITH, S. A.; BEAULIEU, J. M.; DONOGHUE, M. J. Mega-phylogeny approach for comparative biology: An alternative to supertree and supermatrix approaches. **BMC Evolutionary Biology**, v. 9, n. 1, p. 1–12, 2009.

SMITH, S. A.; BROWN, J. W. Constructing a broadly inclusive seed plant phylogeny. **American Journal of Botany**, v. 105, n. 3, p. 1–13, 2018.

SOKAL, R.; MICHENER, C. A statistical method for evaluating systematic relationships. **The University of Kansas science bulletin**, v. 38, n. 2, 1958.

STAMATAKIS, A. RAxML-VI-HPC: Maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. **Bioinformatics**, v. 22, n. 21, p. 2688–2690, 2006.

STEIN, L. GENOME ANNOTATION: FROM SEQUENCE TO BIOLOGY. **NATURE REVIEWS | GENETICS**, v. 2, p. 493–505, 2001. Disponível em: <[http://www.nature.com/nrg/journal/v2/n7/pdf/nrg0701\\_493a.pdf](http://www.nature.com/nrg/journal/v2/n7/pdf/nrg0701_493a.pdf)>. Acesso em: 20/6/2016.

STEPHENS, Z. D.; LEE, S. Y.; FAGHRI, F.; et al. Big data: Astronomical or genomics? **PLoS Biology**, v. 13, n. 7, p. 1–11, 2015.

STEVENS, P. F. Angiosperm Phylogeny Website. Disponível em: <<http://www.mobot.org/MOBOT/research/APweb/>>. Acesso em: 1/3/2018.

SWOFFORD, DA. L. PAUP\*. Phylogenetic Analysis Using Parsimony (\*and Other Methods). **Journal of Molecular Evolution**, 2002. Sinauer Associates, Sunderland, Massachusetts.

THE ANGIOSPERM PHYLOGENY GROUP VI. An update of the Angiosperm Phylogeny Group classification for the orders and families of flowering plants: APG IV. **Botanical Journal of the Linnean Society**, v. 161, n. 2, p. 105–121, 2016.

TROY, C. S.; MACHUGH, D. E.; BAILEY, J. F.; et al. Genetic evidence for Near-Eastern origins of European cattle. **Nature**, v. 410, n. 6832, p. 1088–1091, 2001.

VINGA, S.; ALMEIDA, J. Alignment-free sequence comparison - A review. **Bioinformatics**, 2003.

WATSON, J. D.; CRICK, F. H. C. Molecular structure of nucleic acids. **Nature**, 1953. Disponível em: <<http://www.nature.com/physics/looking-back/crick/%5Cnhttp://www.ncbi.nlm.nih.gov/pubmed/13054692>>. .

WEEDEN, N. F. Genetic and biochemical implications of the endosymbiotic origin of the chloroplast. **Journal of Molecular Evolution**, v. 17, n. 3, p. 133–139, 1981.

WOLFE, K. H.; MORDENT, C. W.; PALMERS, J. D. Function and evolution of a minimal plastid genome from a nonphotosynthetic parasitic plant (chloroplast DNA/Epifagus virginiana/transtdon/rascriptlon/photosyntbesis). **Evolution**, v. 89, p. 10648–10652, 1992.

WOLFSBERG, T. G.; SCHAFER, S.; TATUSOV, R. L.; TATUSOVA, T. A. Organelle genome resources at NCBI. **Trends in Biochemical Sciences**, v. 26, n. 3, p. 199–203, 2001.

YANDELL, M.; ENCE, D. A beginner's guide to eukaryotic genome annotation. **Nature Publishing Group**, v. 13, 2012.

YANG, A.; TROUP, M.; HO, J. W. K. Scalability and Validation of Big Data Bioinformatics Software. **Computational and Structural Biotechnology Journal**, v. 15, p. 379–386, 2017. The Authors. Disponível em: <<https://doi.org/10.1016/j.csbj.2017.07.002>>. .

YANG, E. C.; KIM, K. M.; KIM, S. Y.; et al. Highly conserved mitochondrial genomes among multicellular red algae of the florideophyceae. **Genome Biology and Evolution**, v. 7, n. 8, p. 2394–2406, 2015.

YANG, Z. Phylogenetic analysis using parsimony and likelihood methods. **Journal of Molecular Evolution**, v. 42, n. 2, p. 294–307, 1996. Disponível em: <<http://www.springerlink.com/index/10.1007/BF02198856>>. .

YANG, Z.; RANNALA, B. Bayesian phylogentic inference using DNA sequences: A Markov chain monte carlo method. **Molecular Biology and Evolution**, v. 14, n. 7, p. 717–724, 1997. Disponível em: <<http://mbe.oxfordjournals.org/cgi/content/abstract/14/7/717>>. .

YU, Z.; PENG, S.; YANG, P. The complete chloroplast genome of the southern Chinese pine *Pinus tabuliformis* (Pinales: Pinaceae). **Mitochondrial DNA Part A: DNA Mapping, Sequencing, and Analysis**, v. 28, n. 1, p. 13–14, 2017.

ZHAO, Y.; LU, D.; HAN, R.; WANG, L.; QIN, P. The complete chloroplast genome sequence of the shrubby cinquefoil *Dasiphora fruticosa* (Rosales: Rosaceae). **Conservation Genetics Resources**, v. 0, n. 0, p. 0, 2017. Springer Netherlands. Disponível em: <<http://link.springer.com/10.1007/s12686-017-0899-6>>. .

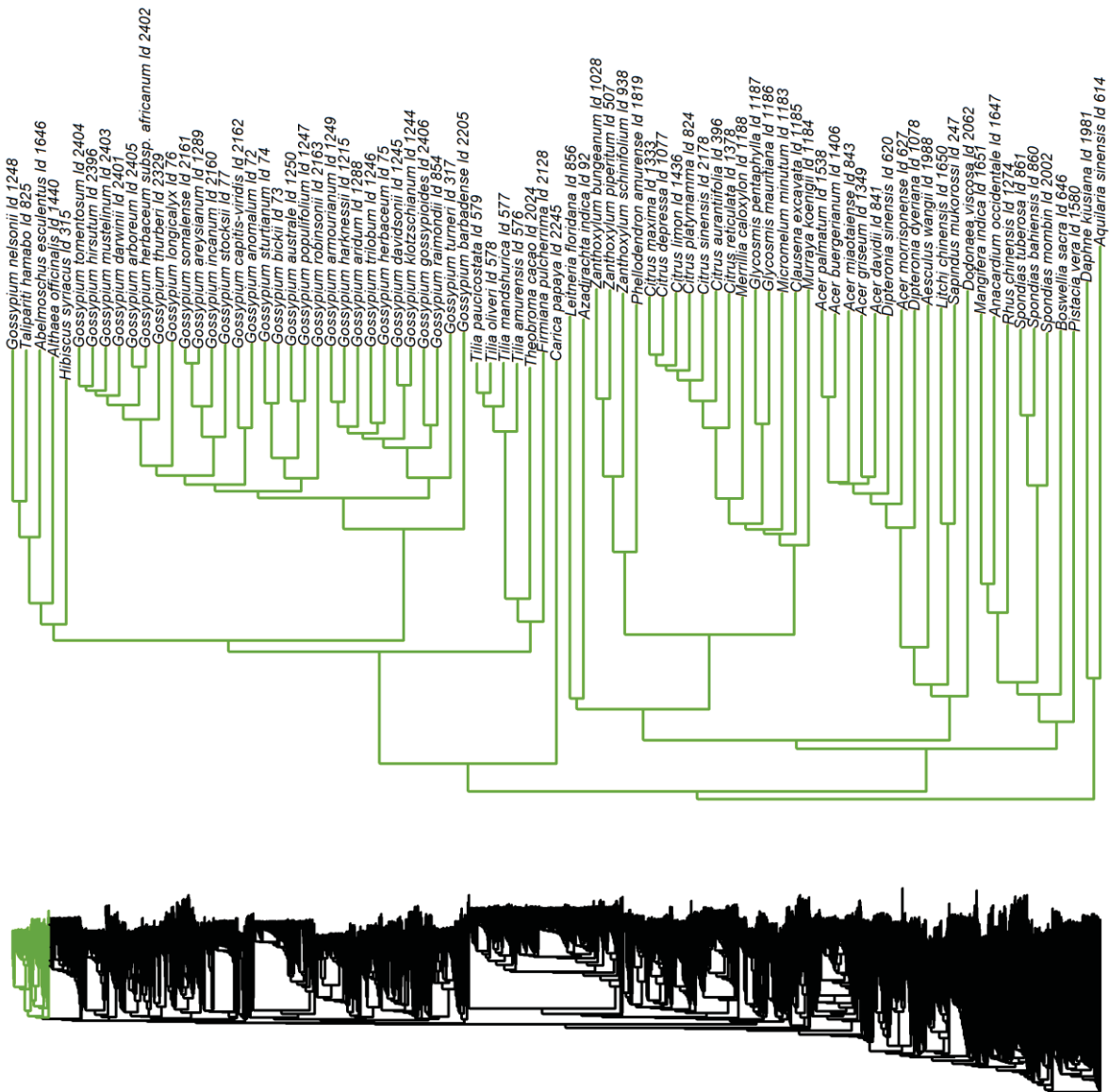
ZUCKERKANDL, E.; PAULING, L. Evolutionary divergence and convergence in proteins. **Evolving Genes and Proteins**, p. 97–166, 1965. Disponível em: <[http://www.yanaiweb.com/genome/Clocks/ZuckerKandl\\_1965.pdf](http://www.yanaiweb.com/genome/Clocks/ZuckerKandl_1965.pdf)>. .

ZVERKOV, O. A.; SELIVERSTOV, A. V.; LYUBETSKY, V. A. PlastidEncoded Protein Families Specific for Narrow Taxonomic Groups of Algae and Protozoa. **Original Russian Text** ©, v. 46, n. 5, p. 717–726, 2012. Pleiades Publishing, Inc.

ZVERKOV, O. A.; SELIVERSTOV, A. V.; LYUBETSKY, V. A. A database of plastid protein families from red algae and apicomplexa and expression regulation of the moeB Gene. **BioMed Research International**, v. 2015, 2015.



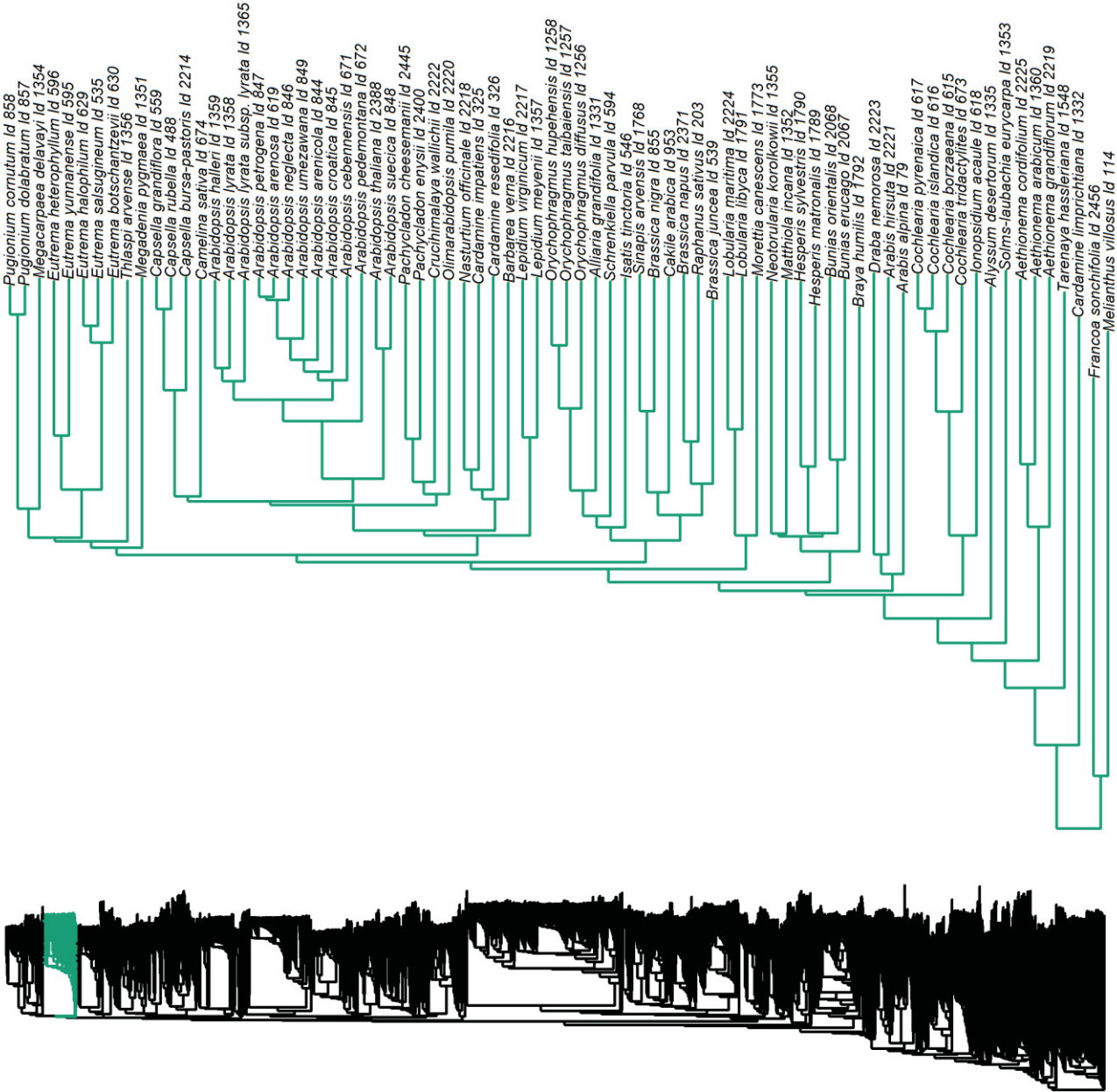
APÉNDICE I : FILOGENIA GLOBAL DE PLASTÍDIOS COMPLETA



Malvales

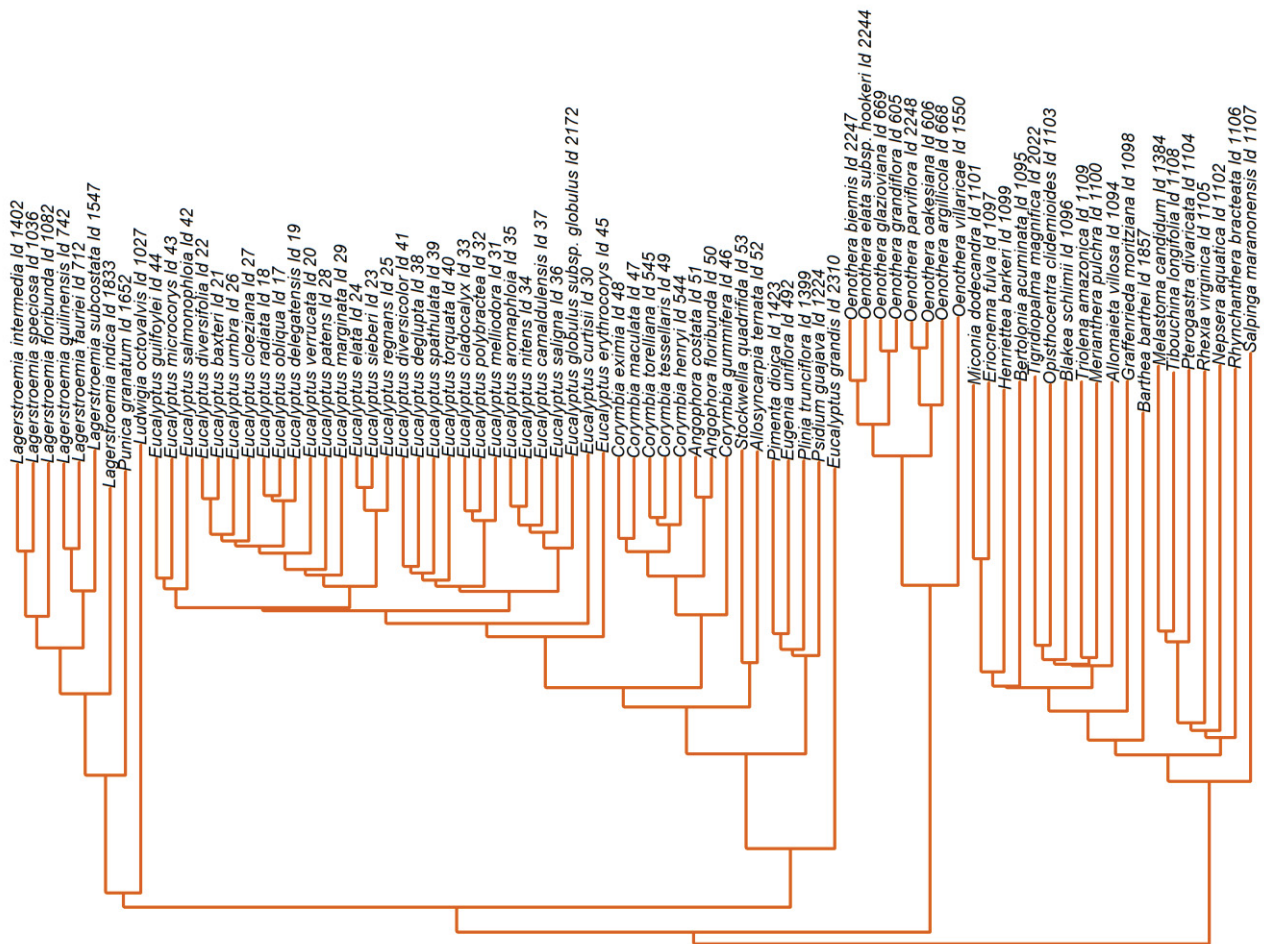
Brassicales

Sapindales



## Brasicales

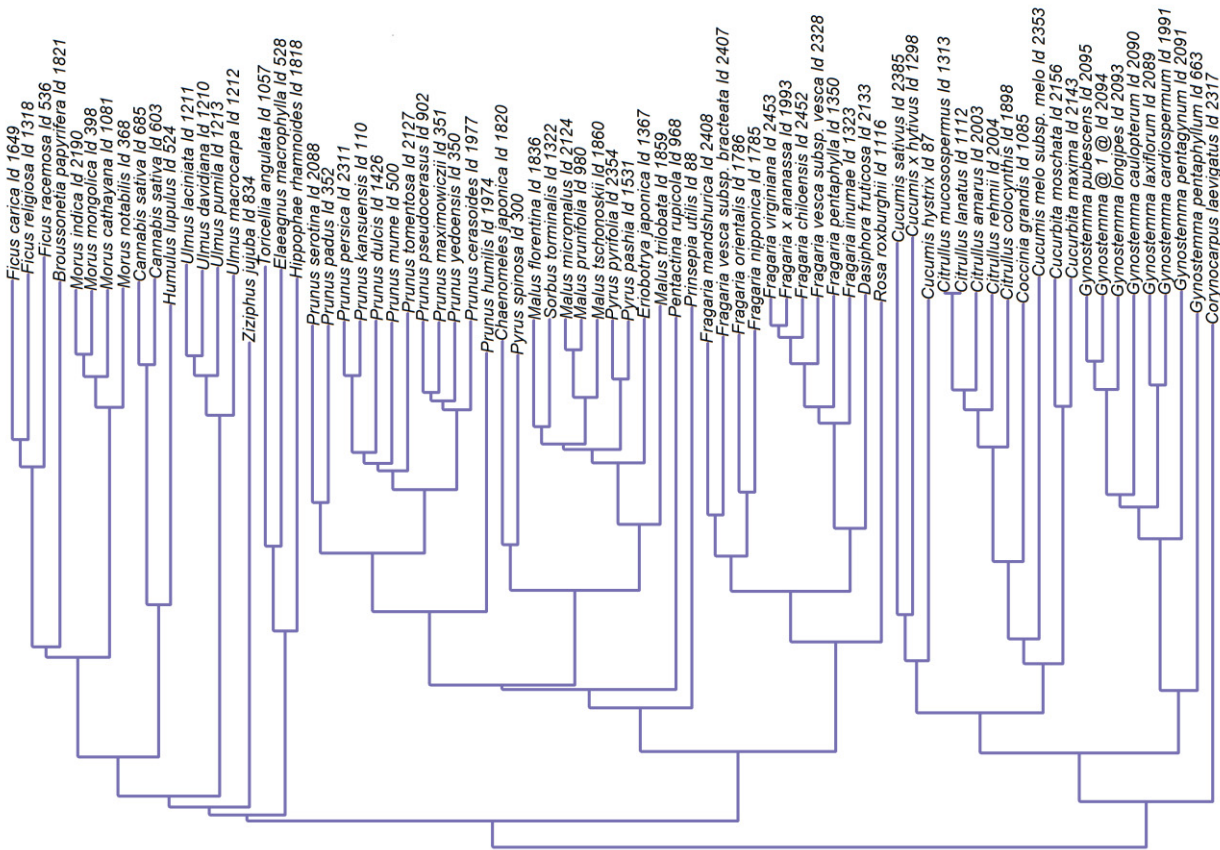
## Geraniales



Myrtales



## Rosales

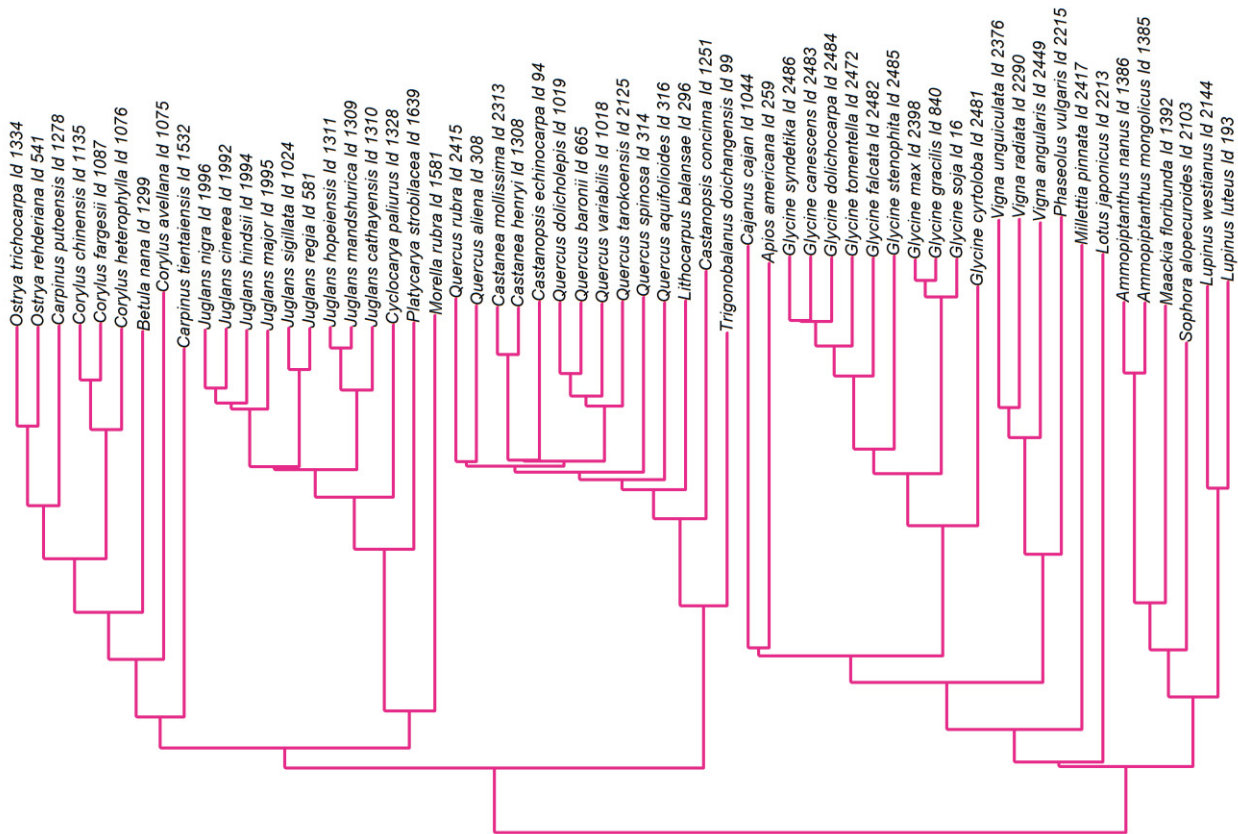


## Cucurbitales

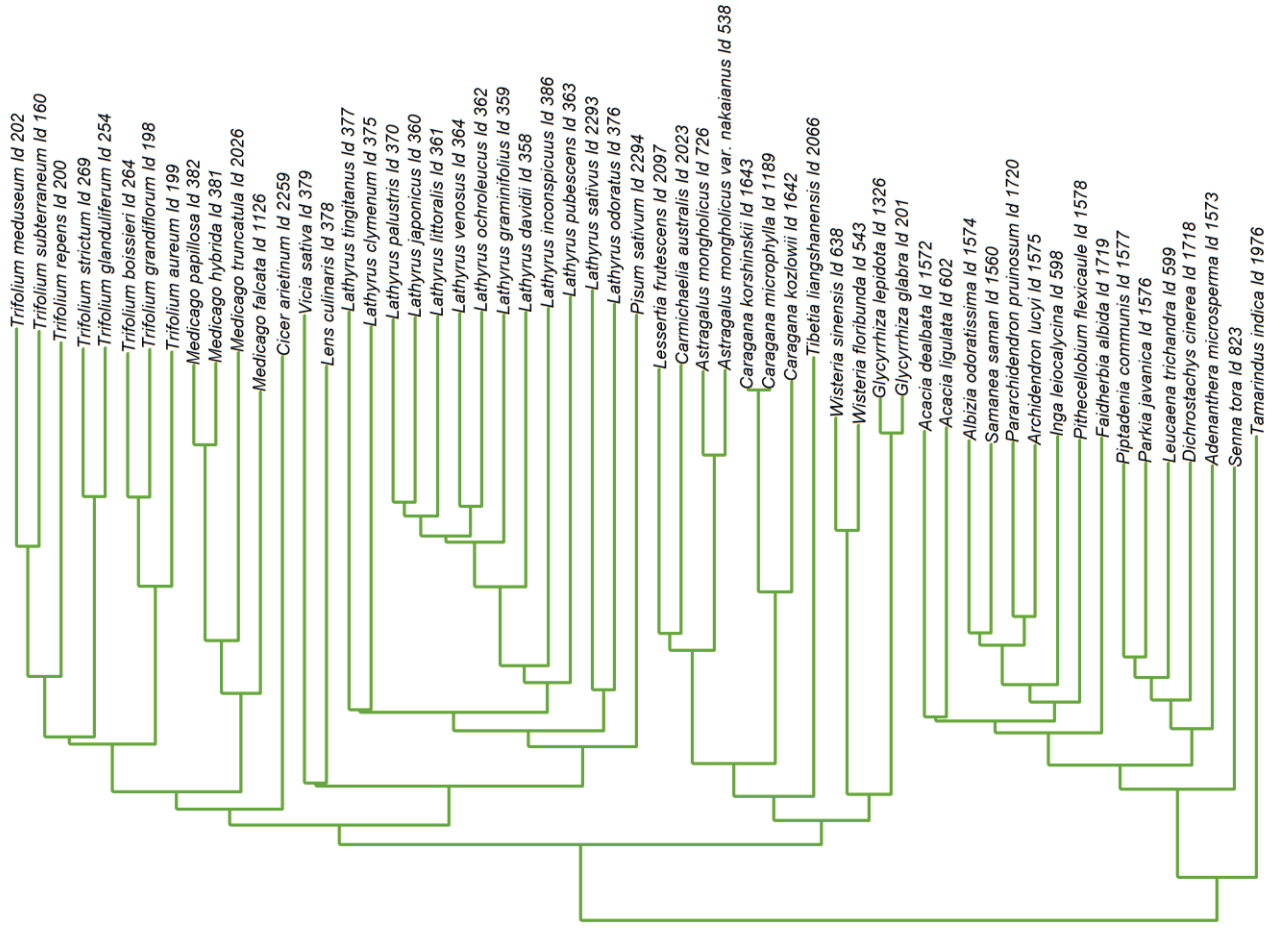




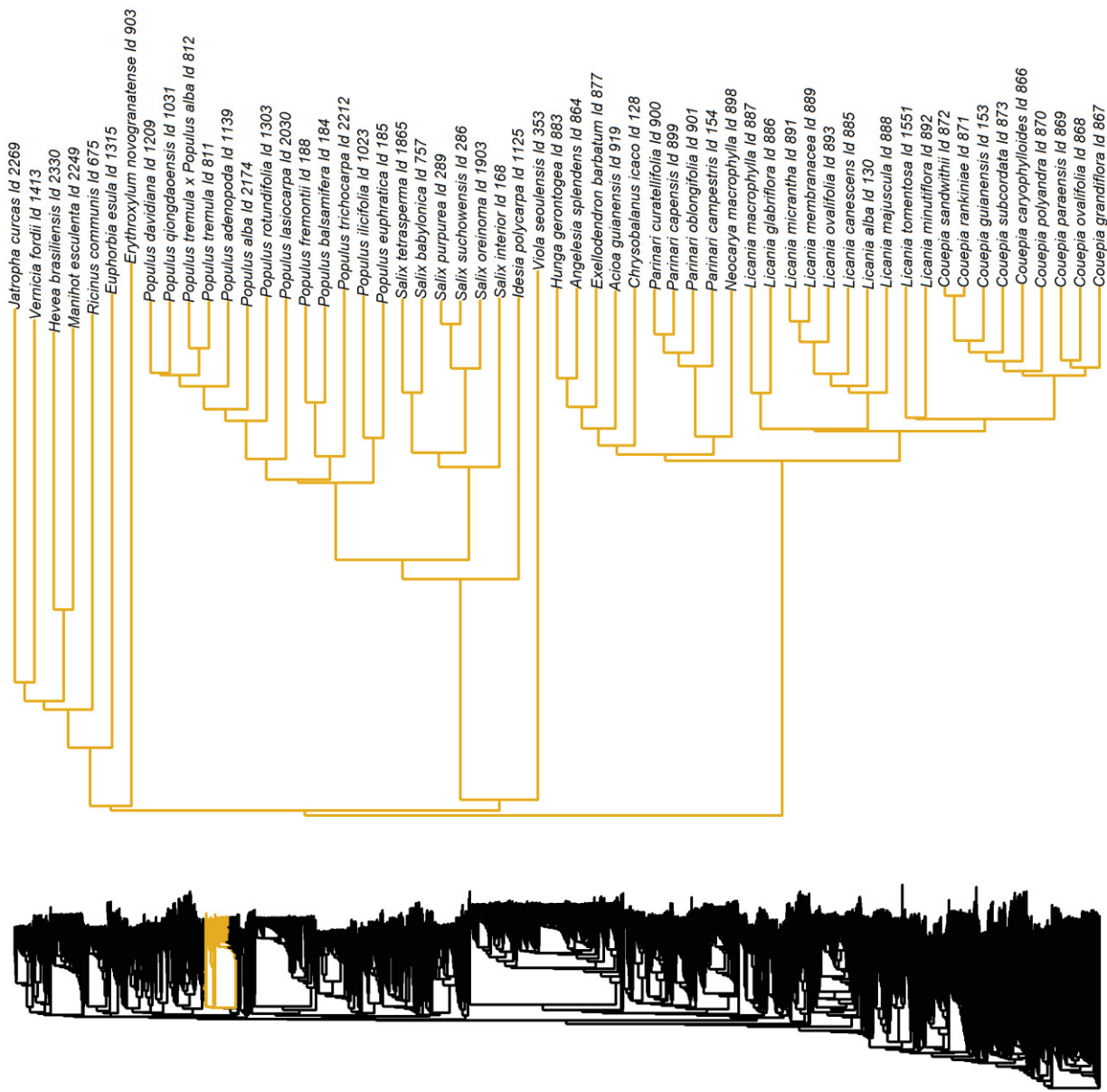
Fagales



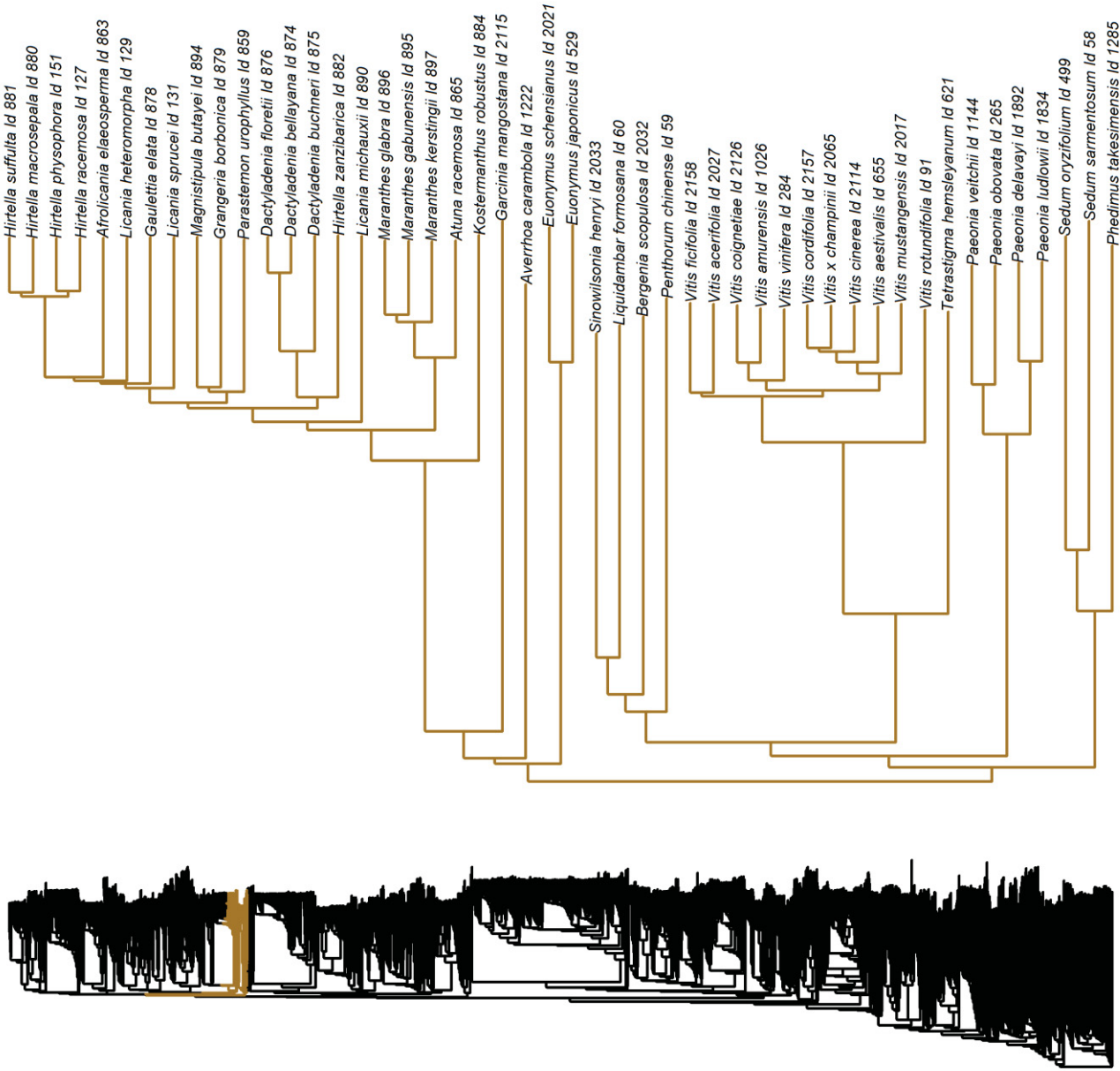
Fabales



## Fabales



## Malpighiales



## Malpighiales

## Oxalidales

## Celastrales

## Saxifragales

## Vitales

## Saxifragales

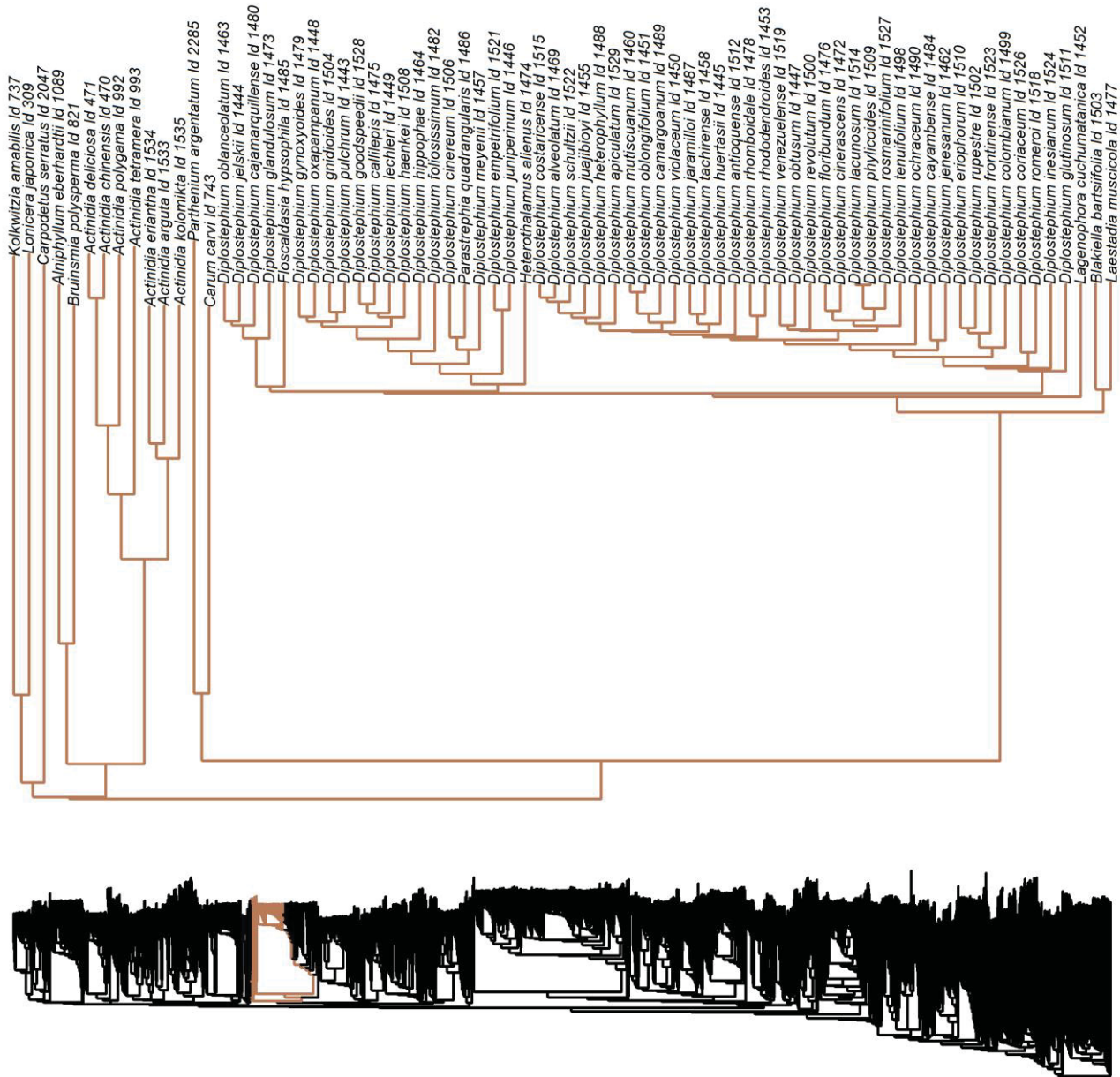


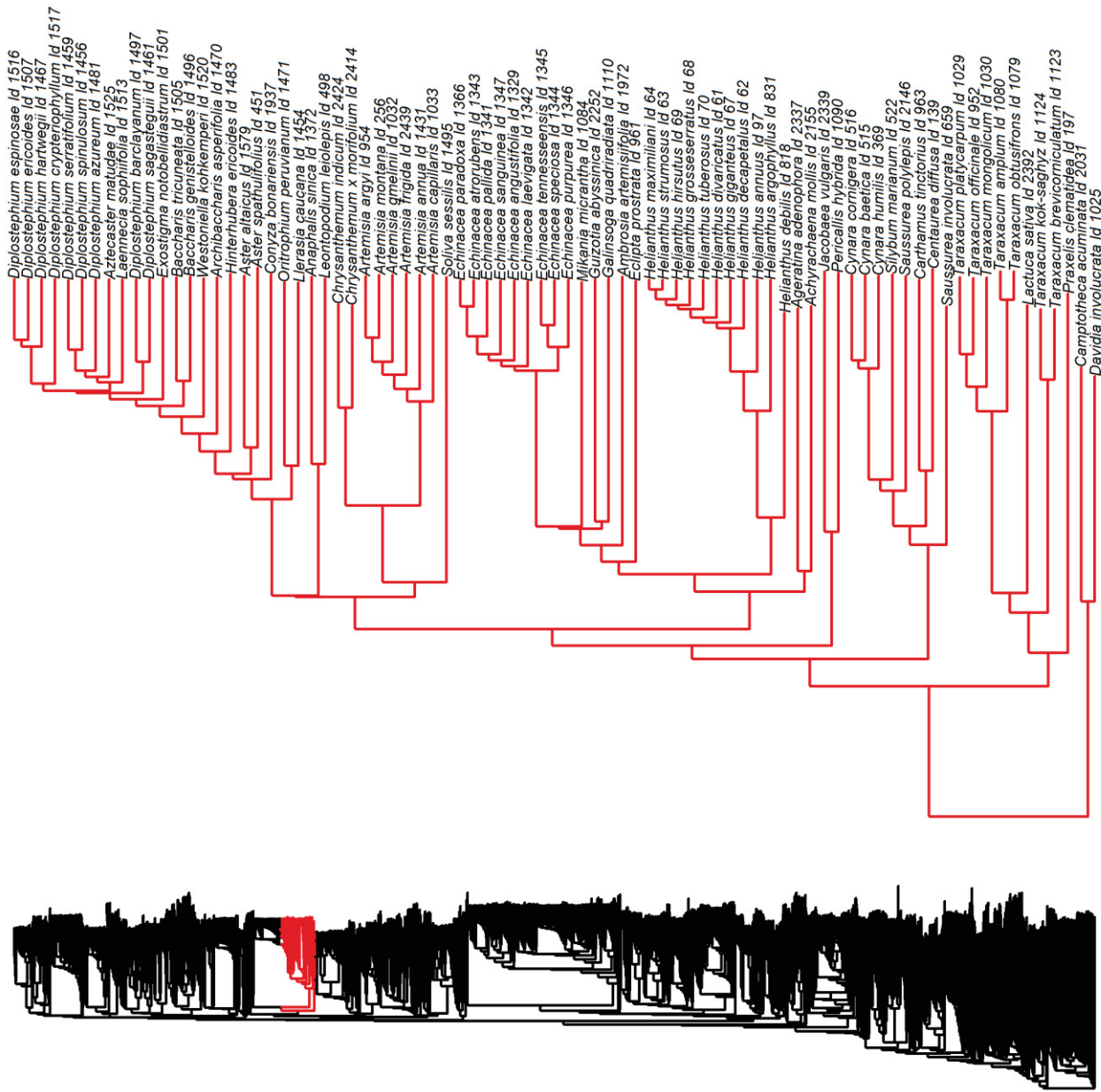
Dipsacales  
Asterales

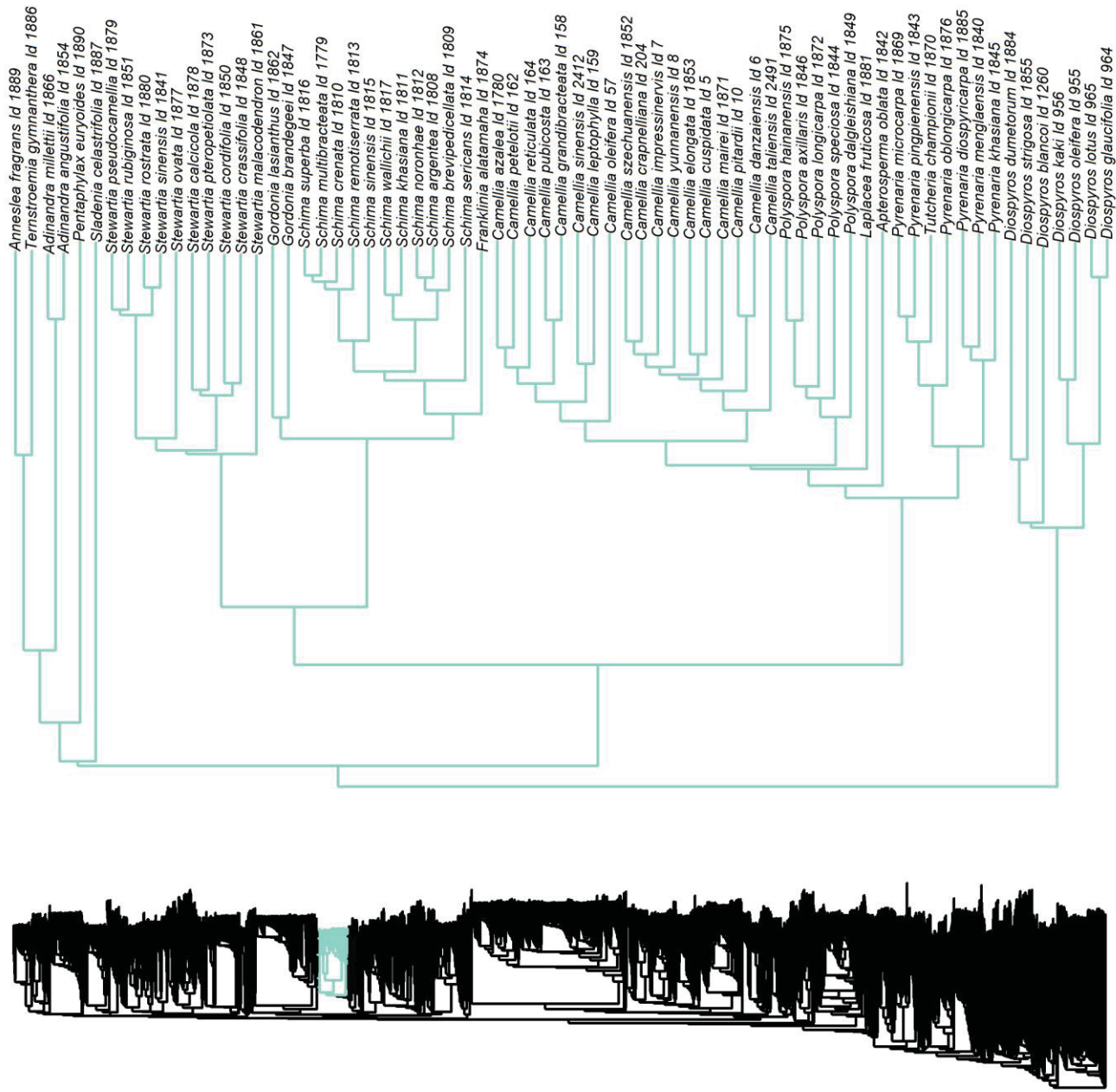
Ericales

Asterales  
Apiales

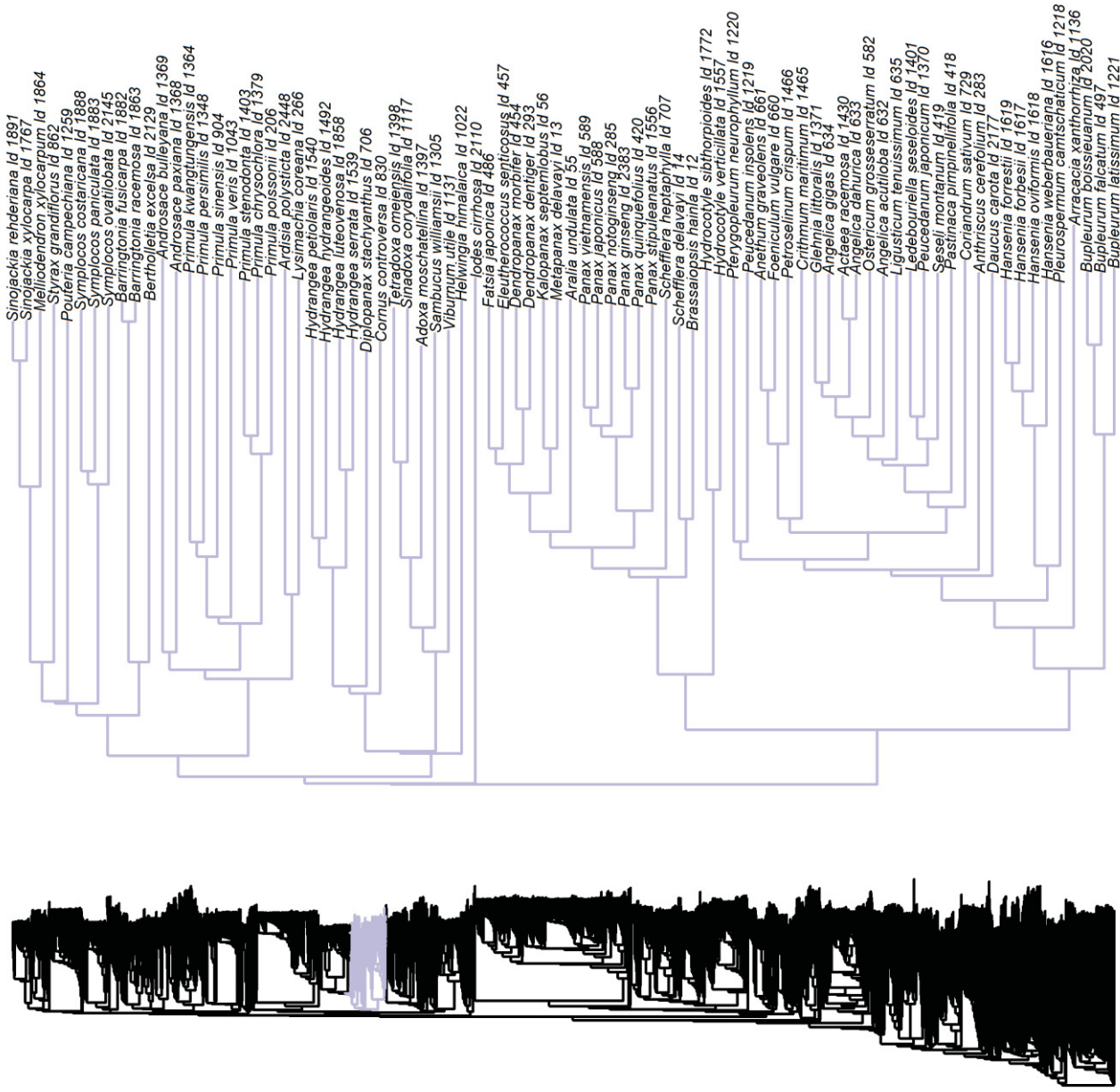
Asterales







## Ericales



## Ericales

## Cornales

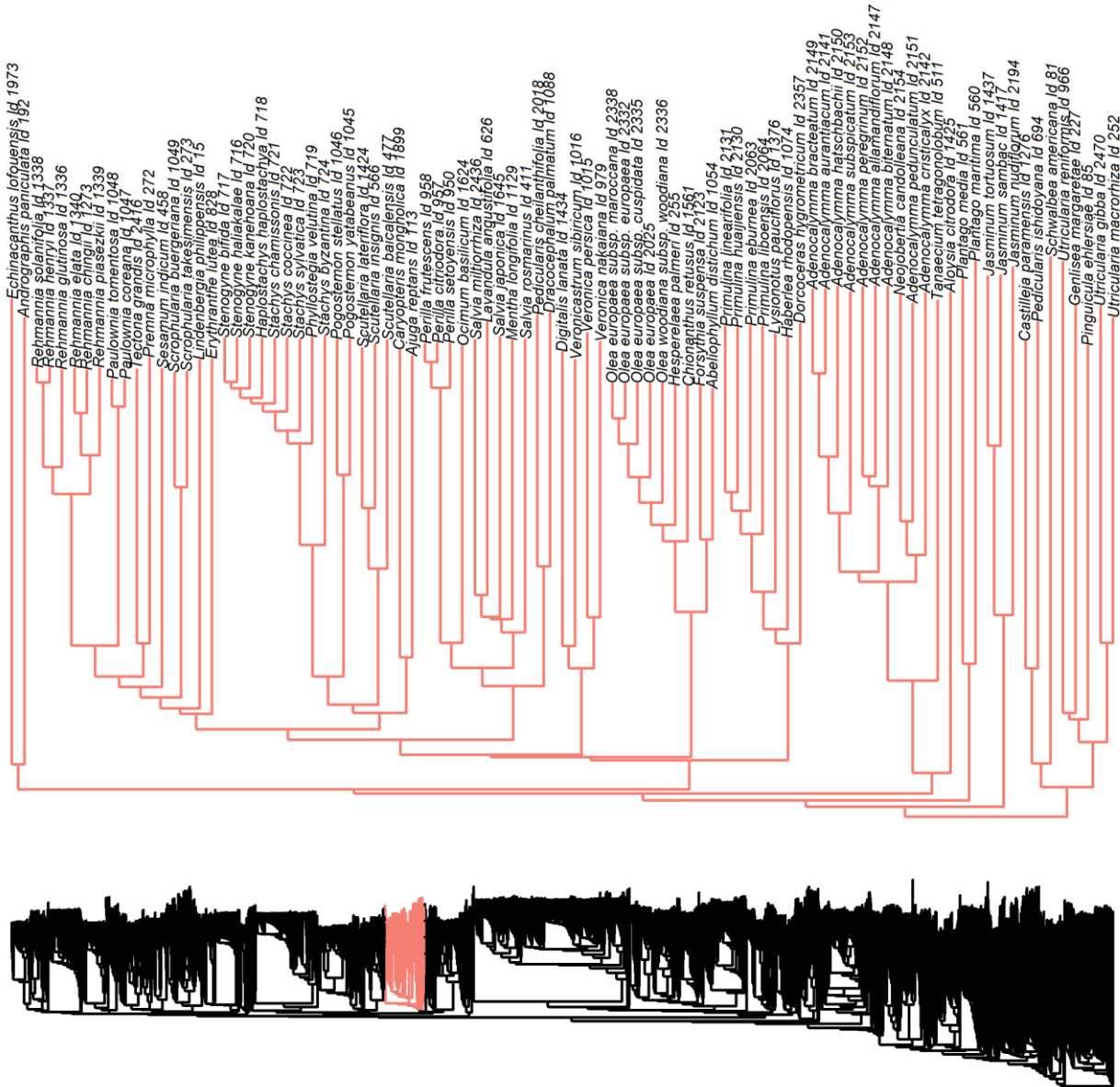
## Dipsacales

## Aquifoliales

## Icacinales

## Apiales



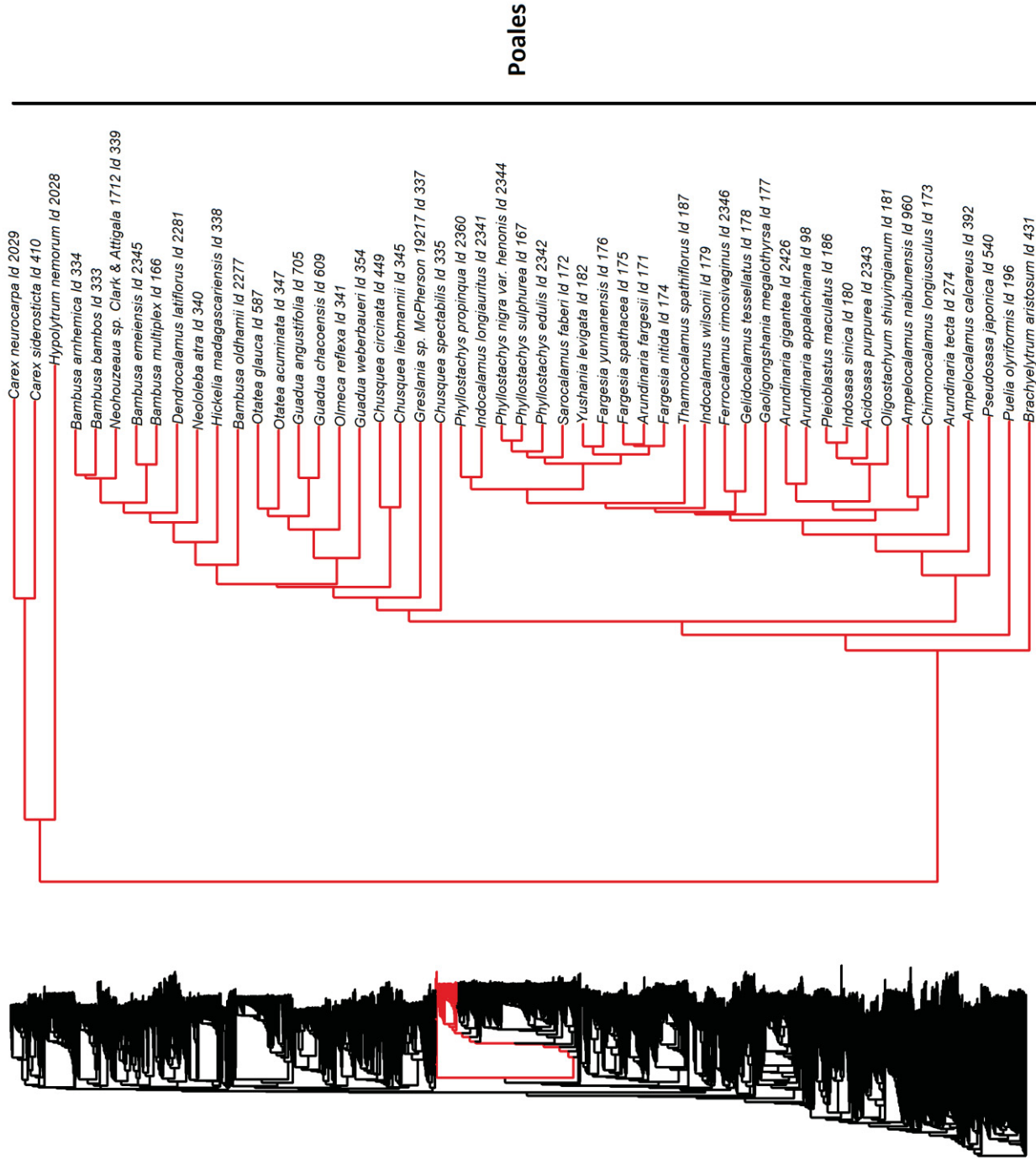


## Lamiales







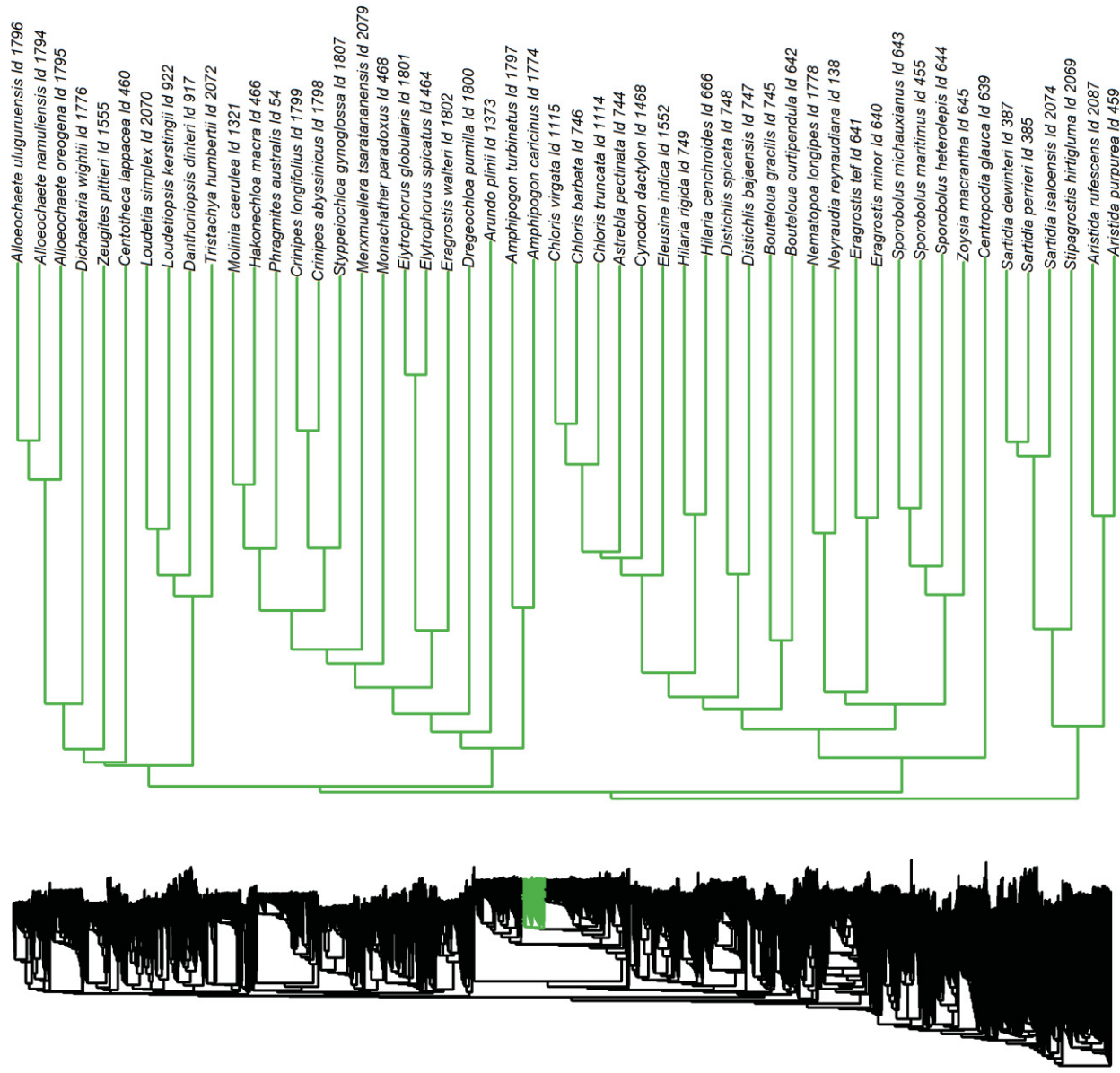


## Poales

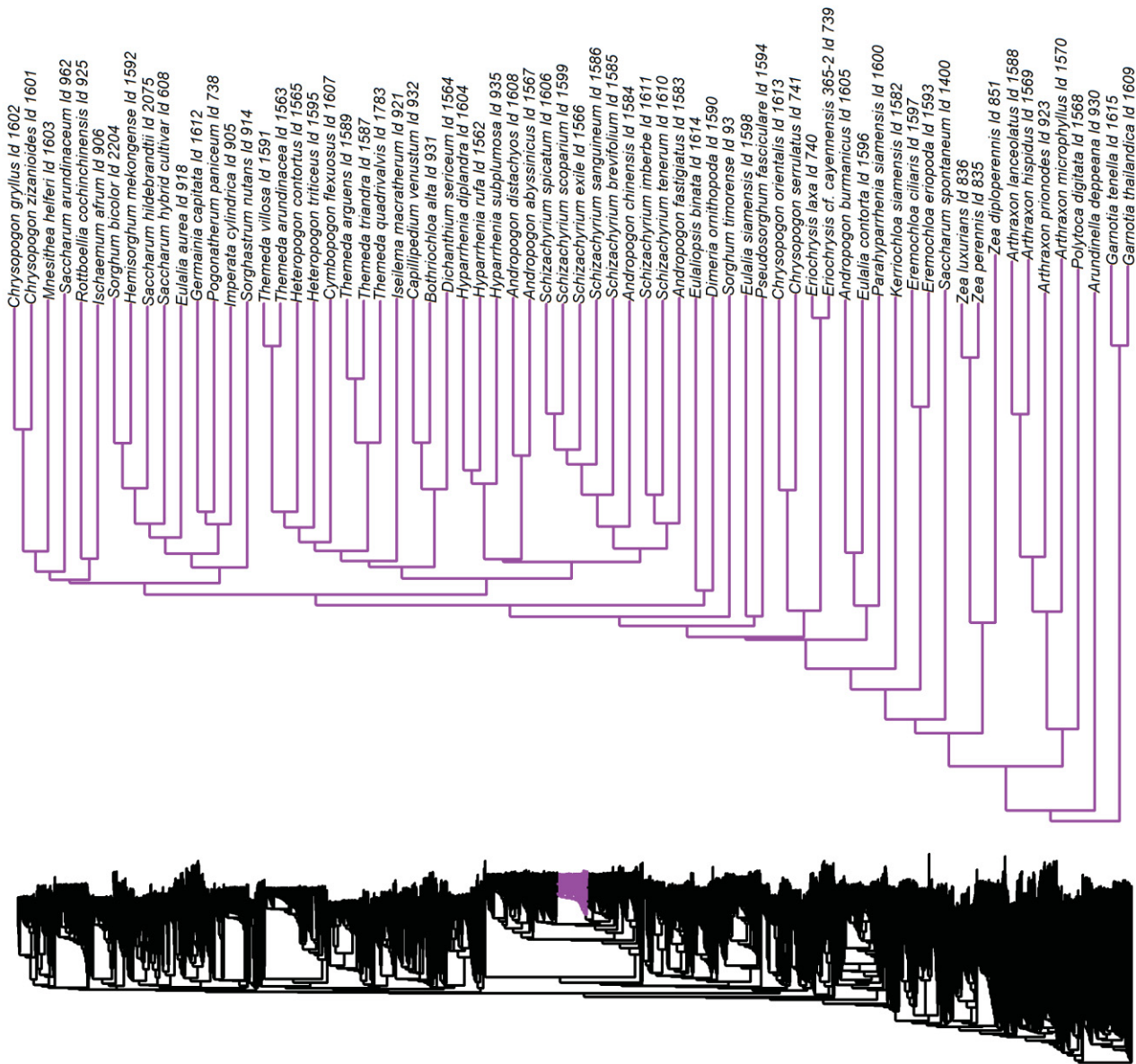


*Torreyochloa pallida* Id 447  
*Phalaris arundinacea* Id 442  
*Agrostis stolonifera* Id 2202  
*Ammophila breviligulata* Id 426  
*Anthoxanthum nitens* Id 436  
*Avena sterilis* Id 1058  
*Avena sativa* Id 429  
*Trisetum cernuum* Id 448  
*Briza maxima* Id 432  
*Anthoxanthum odoratum* Id 428  
*Poa palustris* Id 445  
*Phleum alpinum* Id 443  
*Puccinellia nuttalliana* Id 446  
*Dactylis glomerata* Id 434  
*Festuca ovina* Id 2442  
*Deschampsia antarctica* Id 122  
*Lolium multiflorum* Id 2444  
*Lolium perenne* Id 2239  
*Festuca pratensis* Id 2443  
*Festuca altissima* Id 2441  
*Helictochloa hookeri* Id 430  
*Oryzopsis asperifolia* Id 440  
*Ampelodesmos mauritanicus* Id 427  
*Stipa purpurea* Id 631  
*Piptochaetium avenaceum* Id 444  
*Stipa hymenoides* Id 425  
*Diarrhena obovata* Id 435  
*Phaenosperma globosum* Id 441  
*Melica subulata* Id 439  
*Melica mutica* Id 438  
*Triticum turgidum* Id 502  
*Triticum timopheevii* Id 501  
*Aegilops kotschyi* Id 506  
*Aegilops bicornis* Id 505  
*Aegilops sharonensis* Id 504  
*Aegilops searsii* Id 503  
*Aegilops longissima* Id 494  
*Triticum aestivum* Id 2226  
*Triticum macha* Id 260  
*Triticum monococcum* Id 125  
*Aegilops tauschii* Id 105  
*Aegilops cylindrica* Id 111  
*Aegilops geniculata* Id 112  
*Secale cereale* Id 124  
*Aegilops speltoides* Id 106  
*Littledalea racemosa* Id 2117  
*Bromus vulgaris* Id 433  
*Hordeum vulgare subsp. vulgare* Id 2203  
*Hordeum jubatum* Id 437  
*Triticum urartu* Id 126  
*Brachypodium distachyon* Id 2257  
*Stipa lipskyi* Id 550  
*Diandroyra sp. Clark* Id 1301  
*Olyra latifolia* Id 132  
*Buergerschloa bambusoides* Id 344  
*Pariana campestris* Id 450  
*Pariana radiciflora* Id 348  
*Raddia brasiliensis* Id 342  
*Lithacne pauciflora* Id 346  
*Fargesia denudata* Id 1387

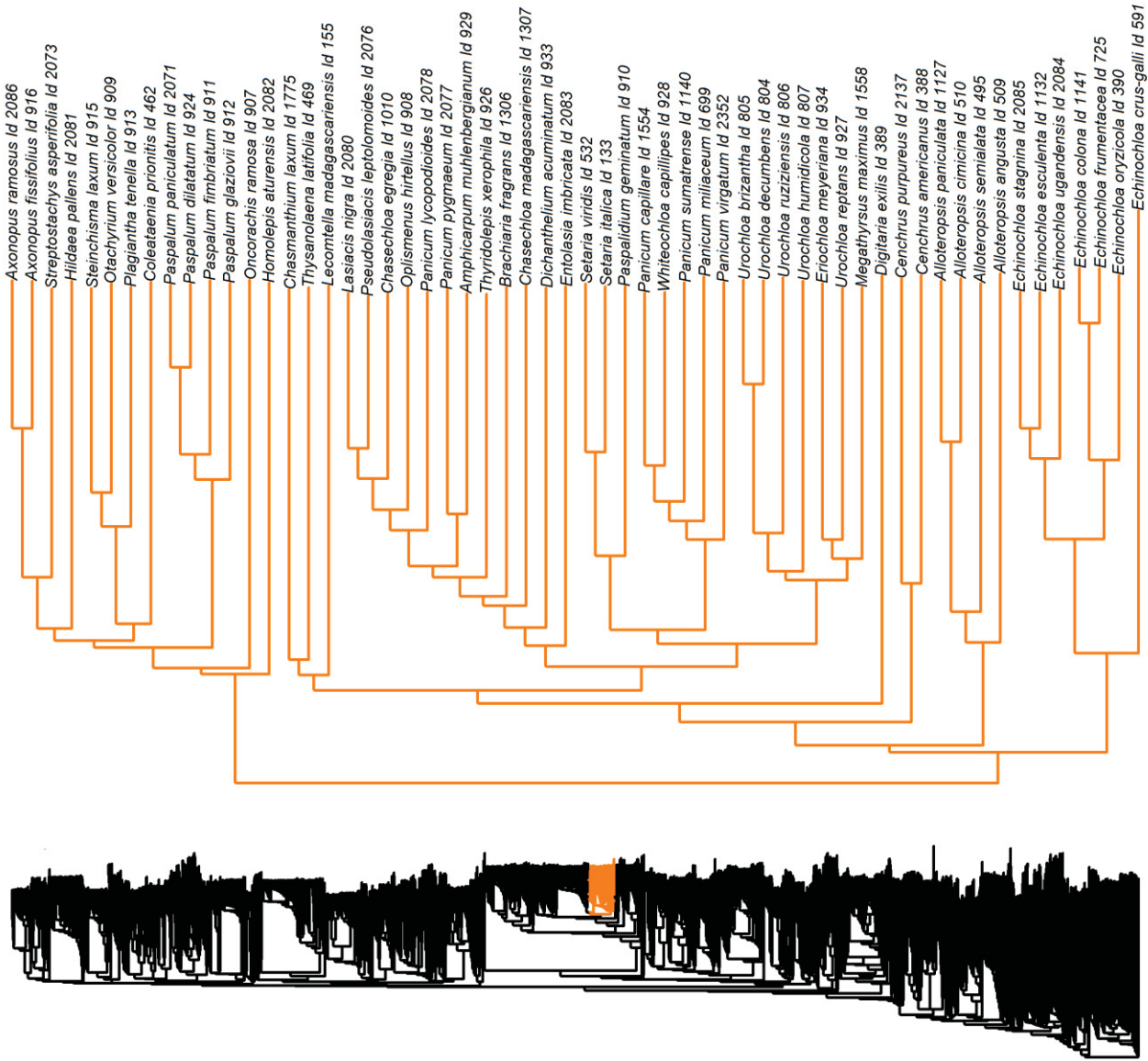
## Poales



## Poales

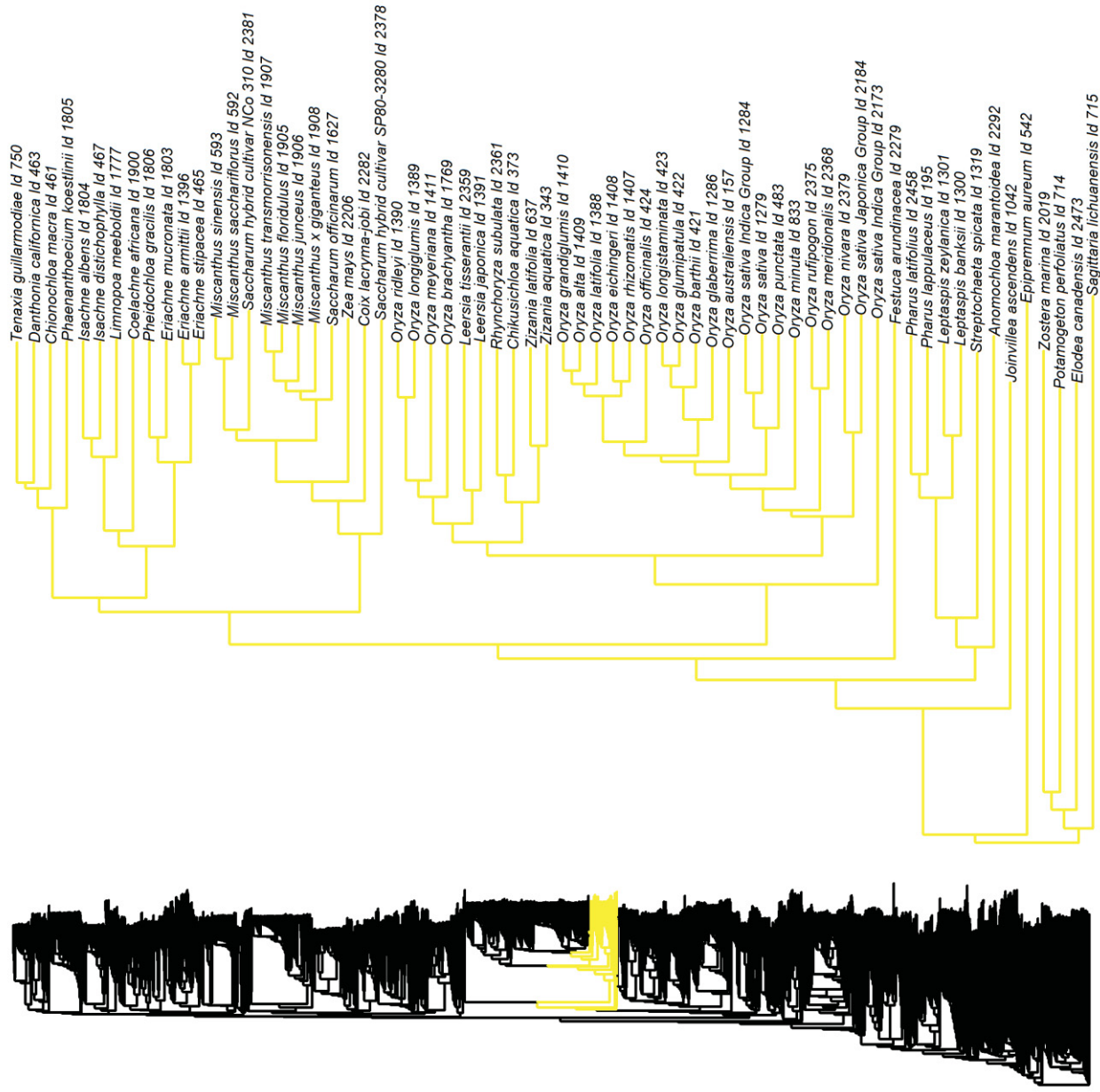


## Poales



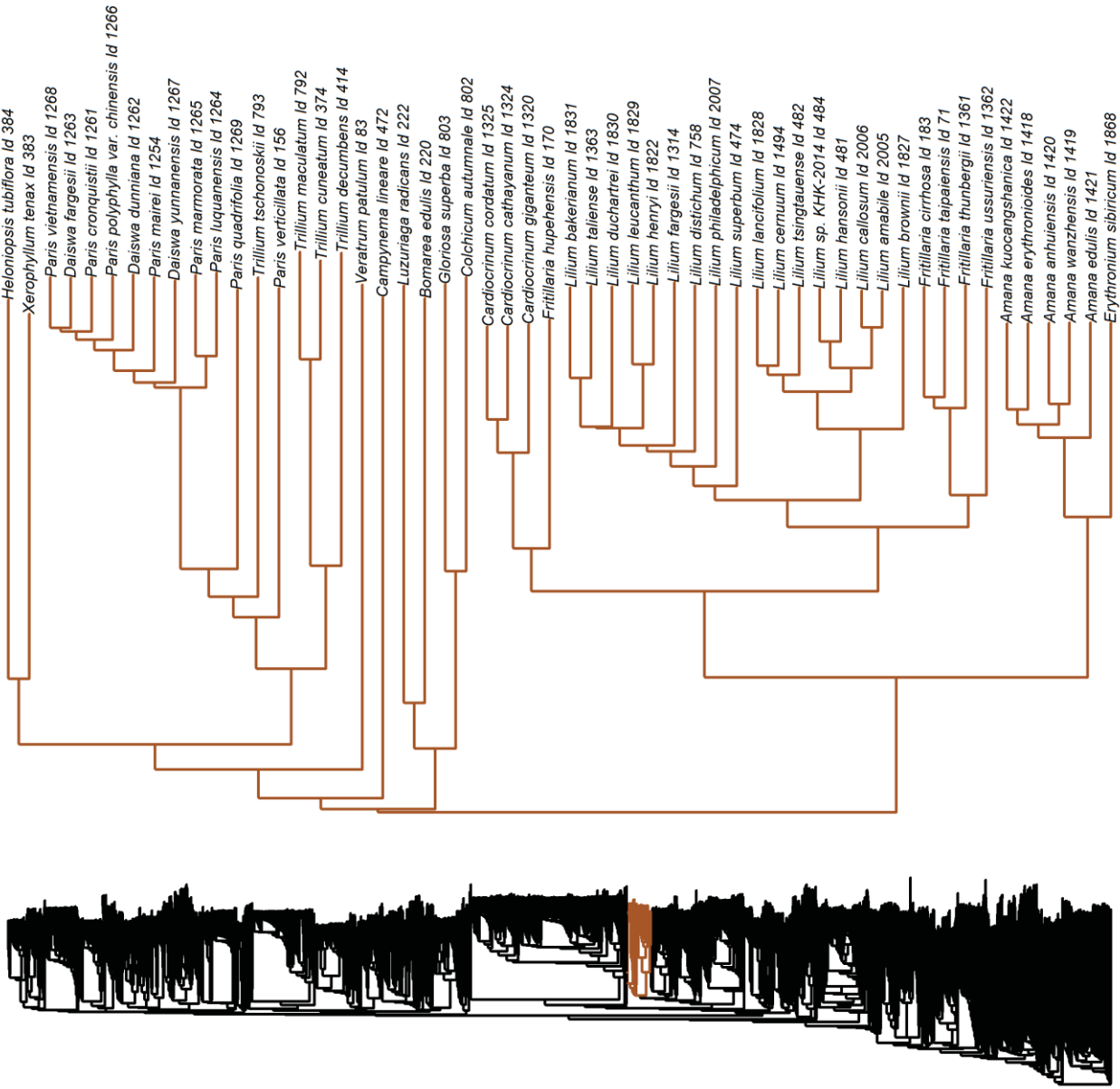
Poales





## Poales





## Liliales

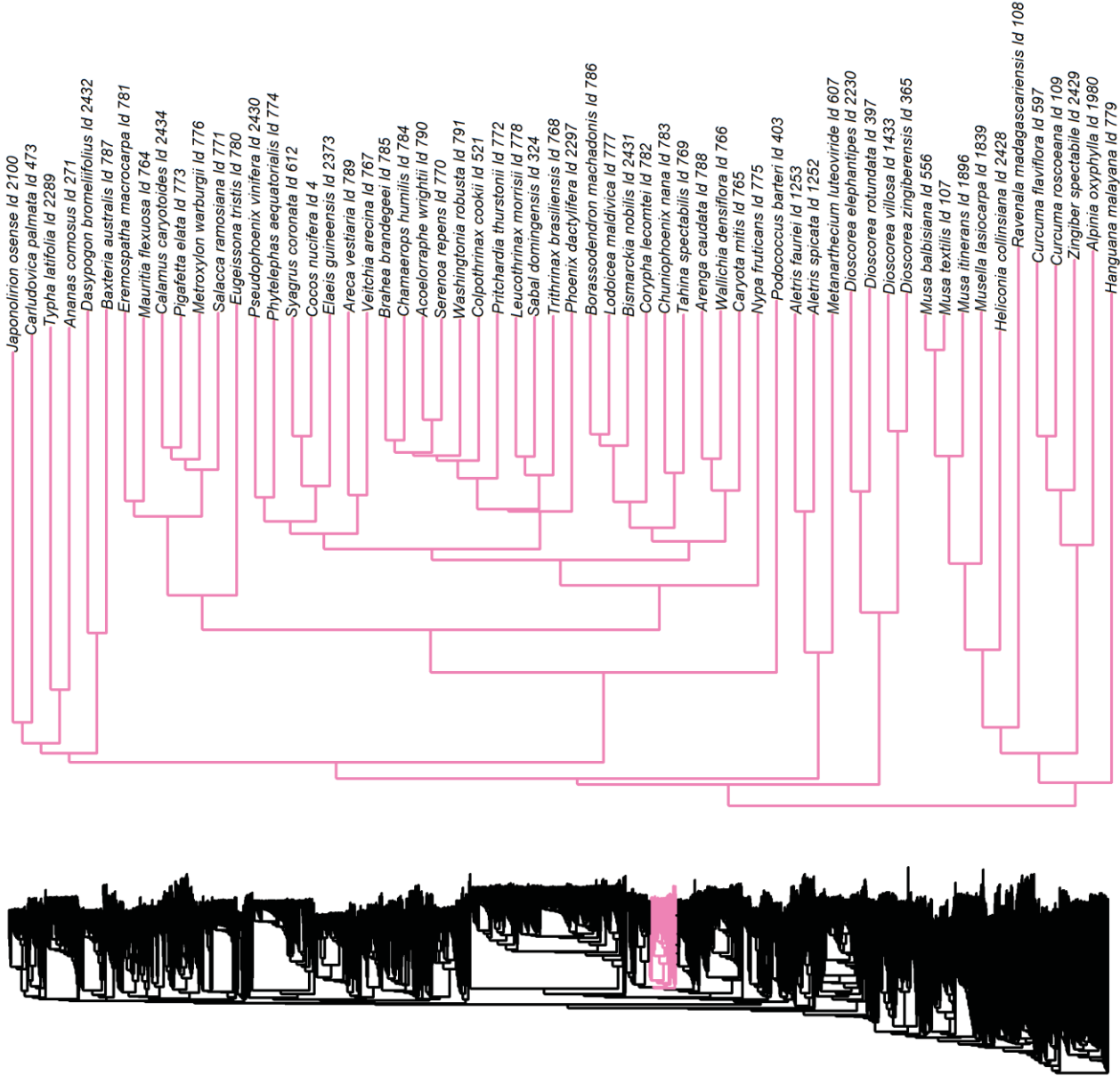
Petrosaviales  
Pandanales  
Poales

Arecales

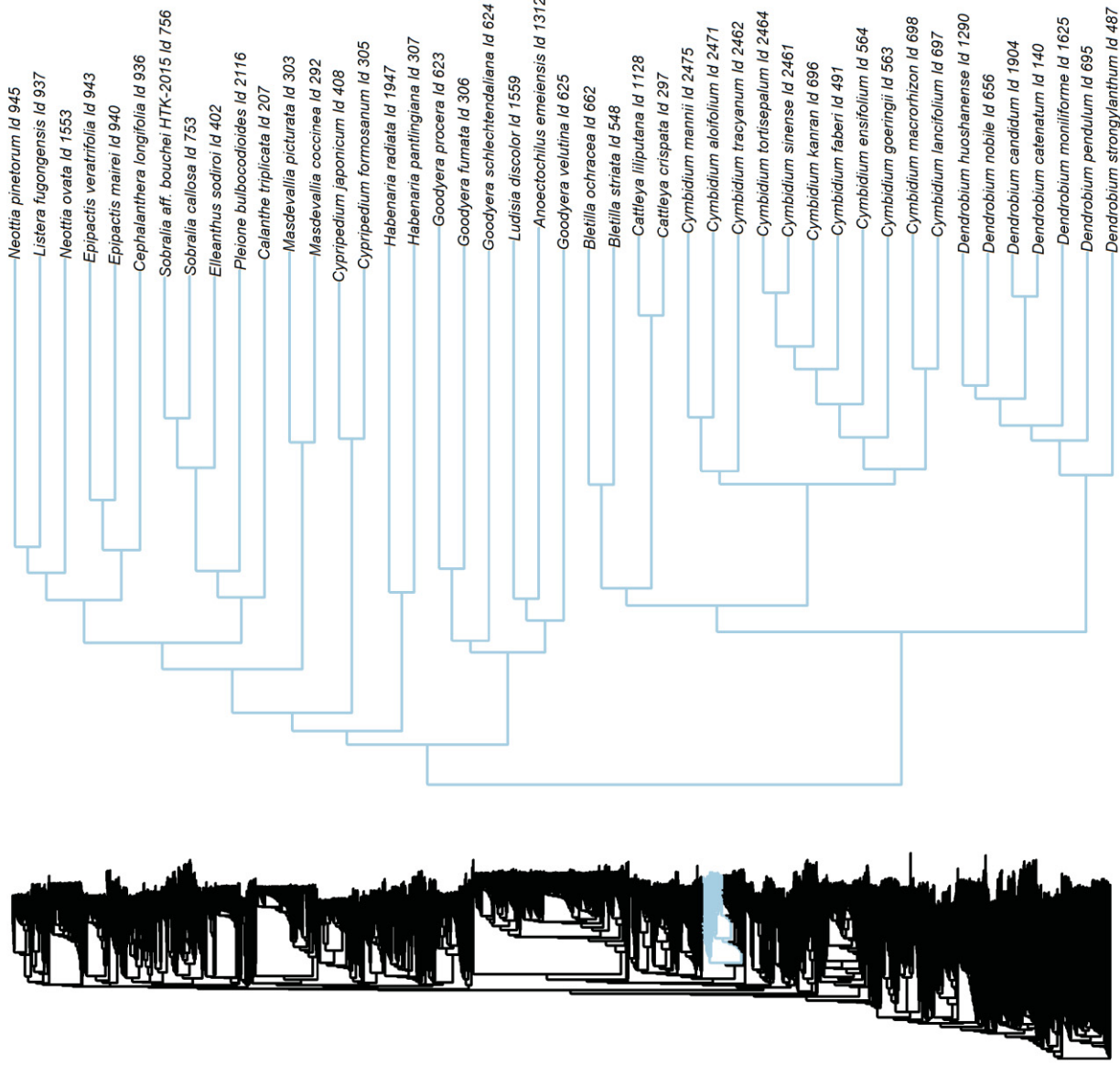
Dioscoreales

Zingiberales

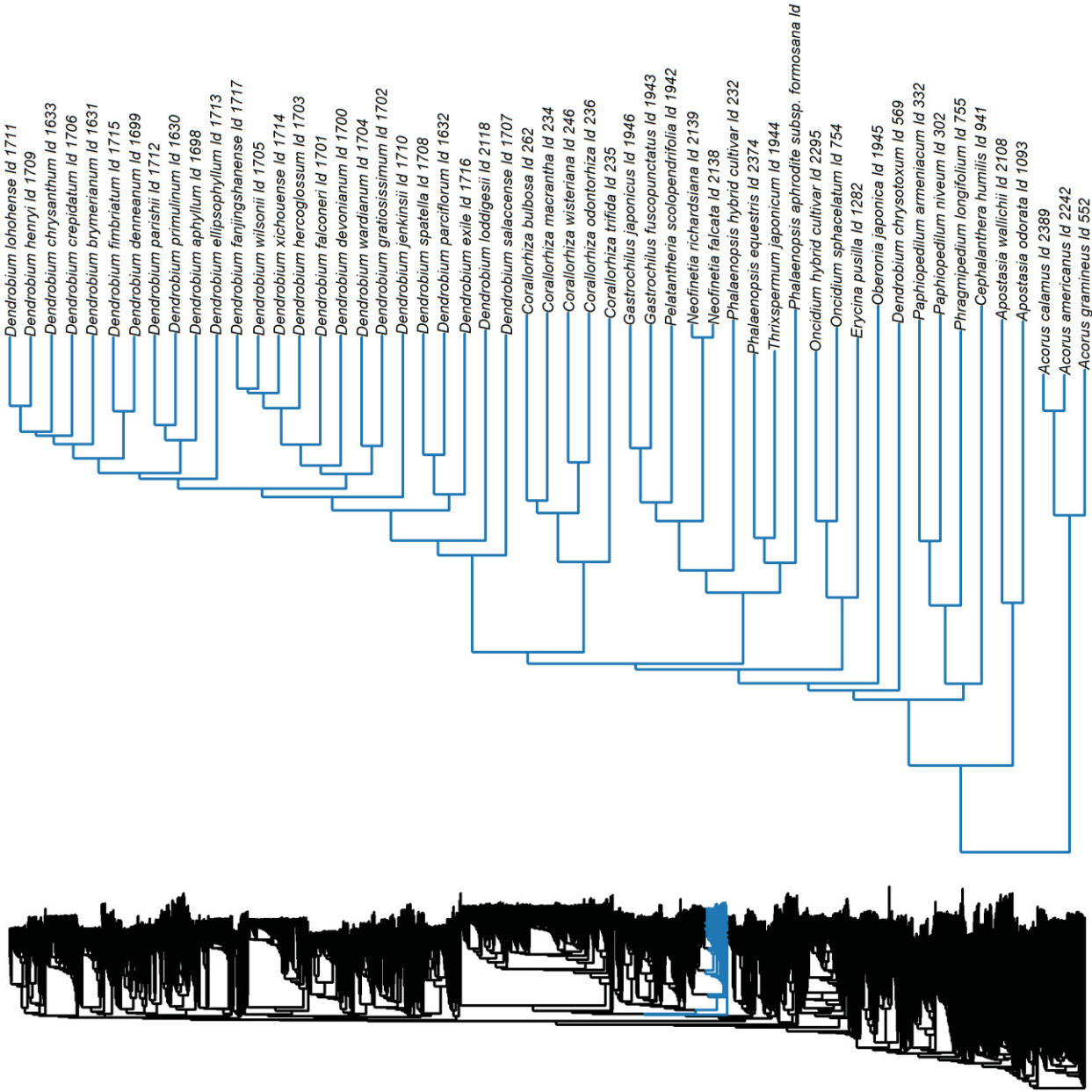
Commelinales





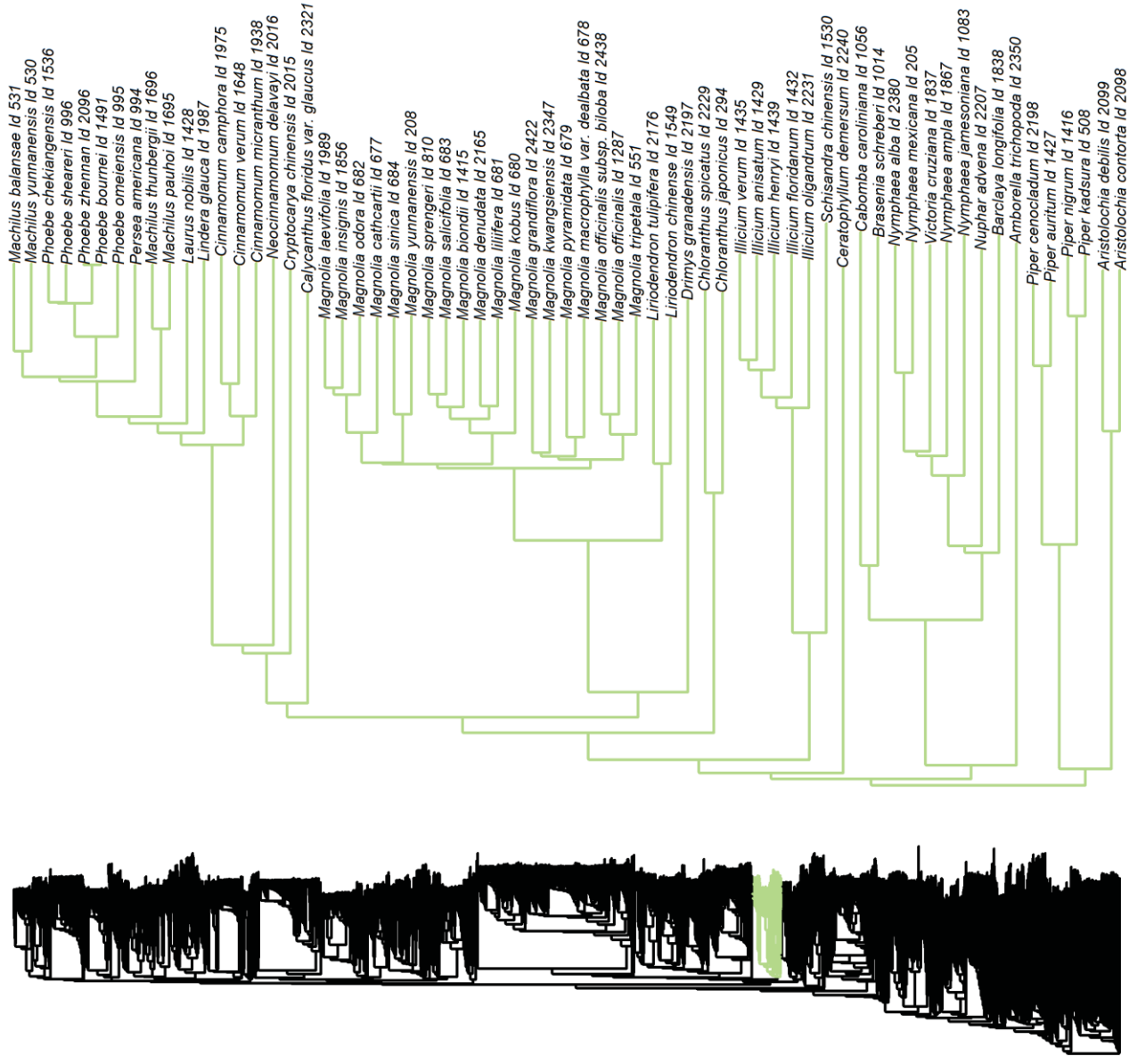


## Asparagales



## Asparagales

## Acorales



## Laurales

## Magnoliales

## Canellales

## Chloranthales

## Austrobaileales

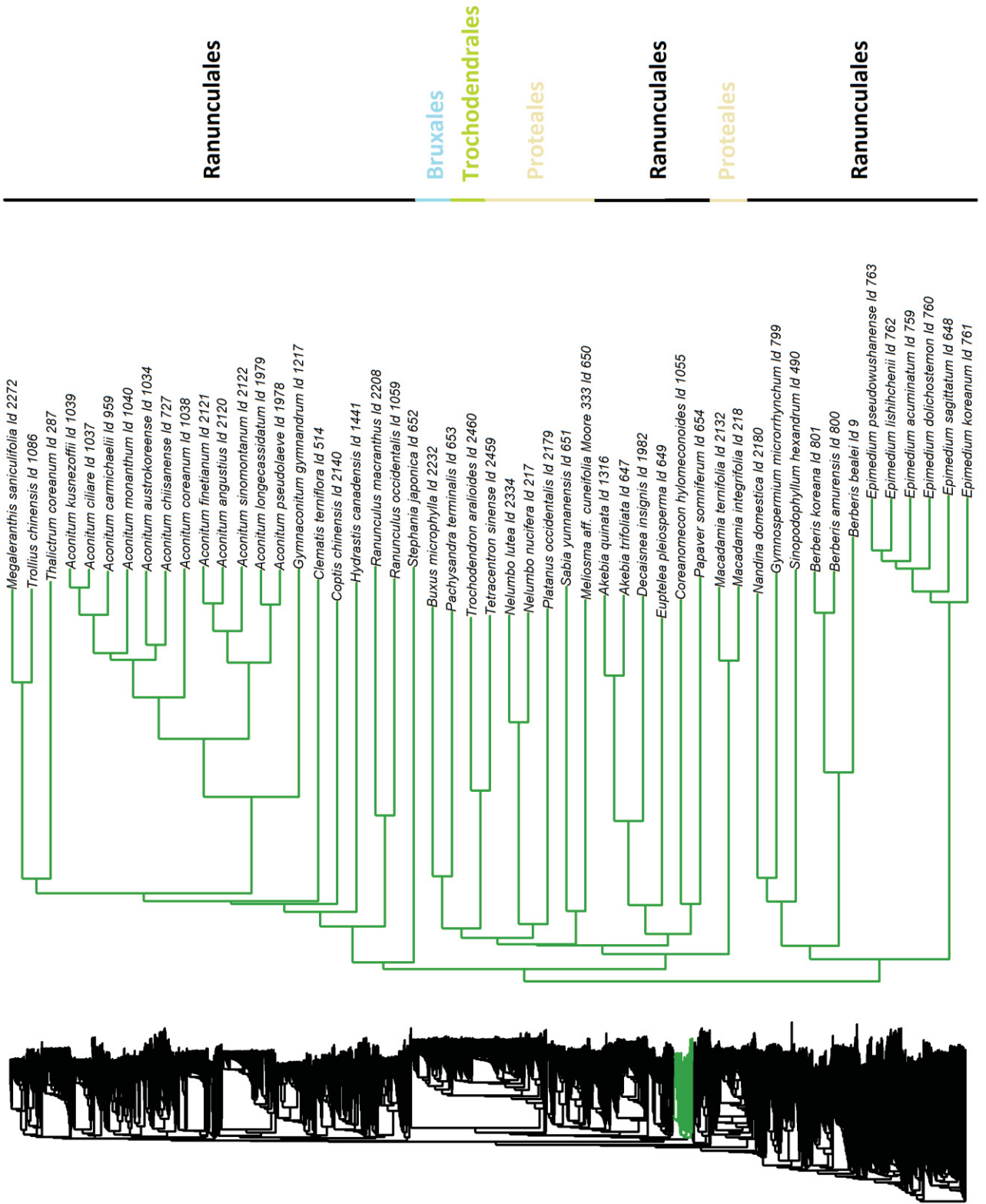
## Ceratophyllales

## Nymphaeales

## Amborellales

## Piperales

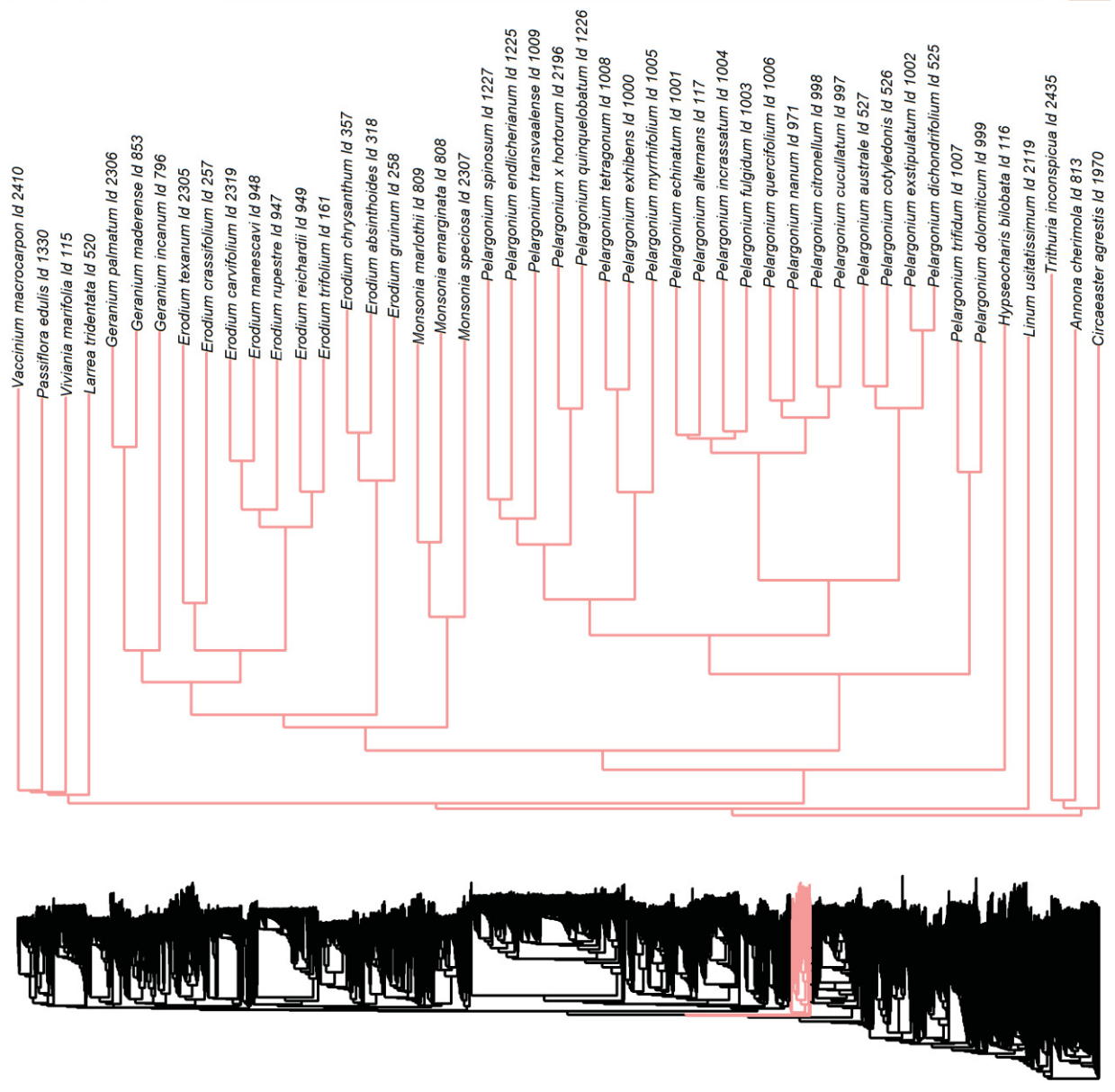


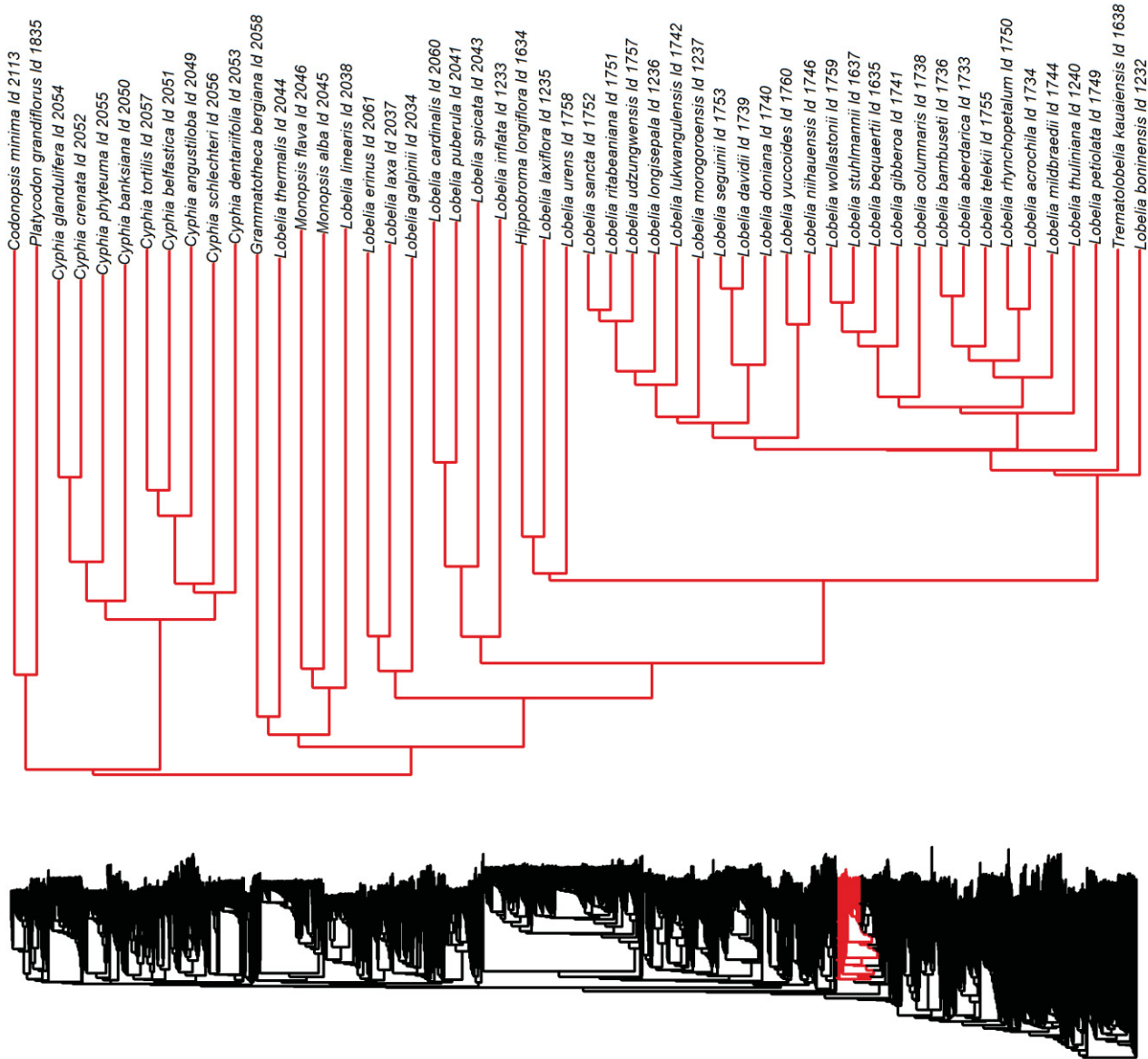


Ericales  
Malpighiales  
Geraniales  
Zygophyllales

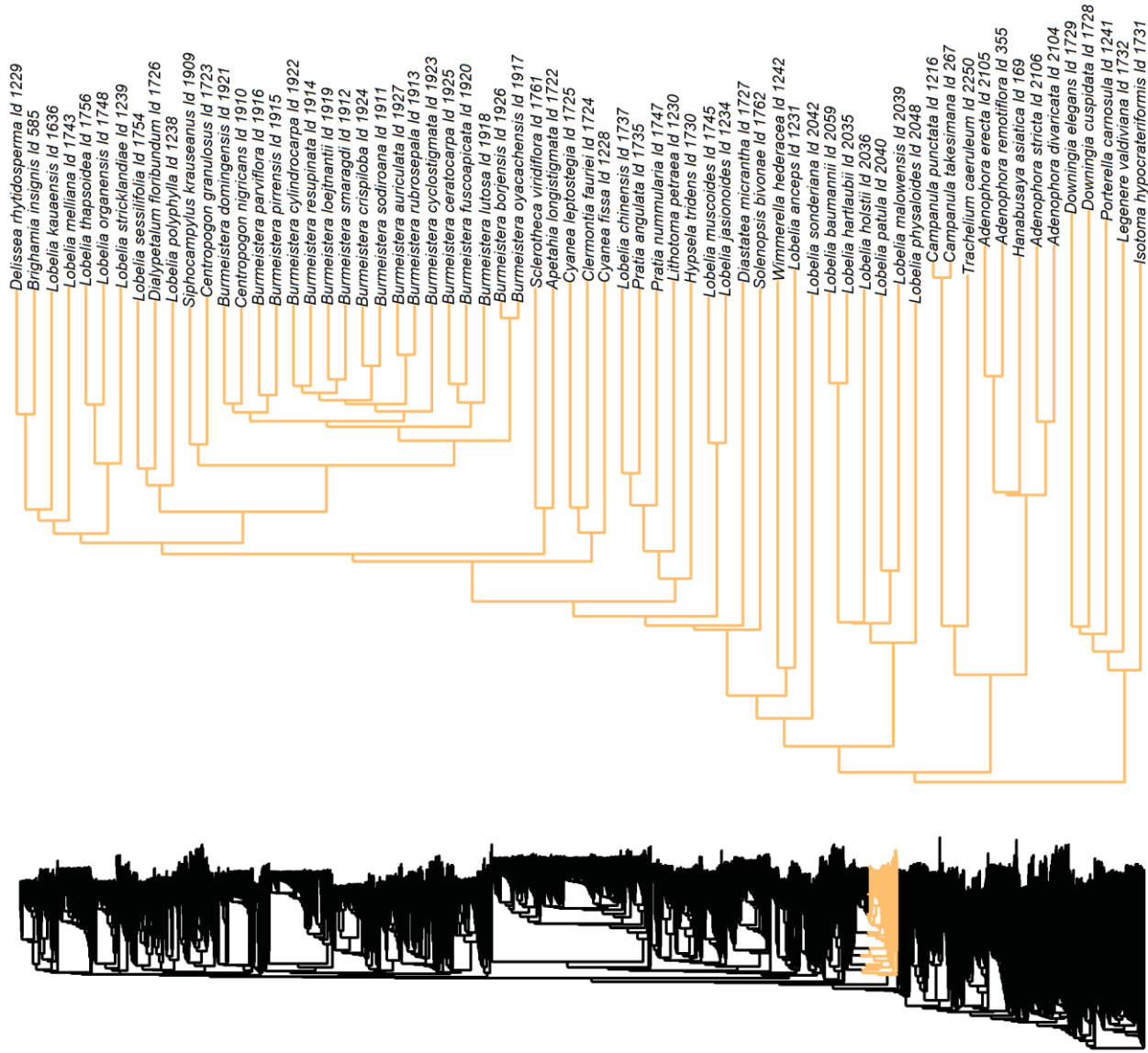
Geraniales

Nymphaeales  
Geraniales

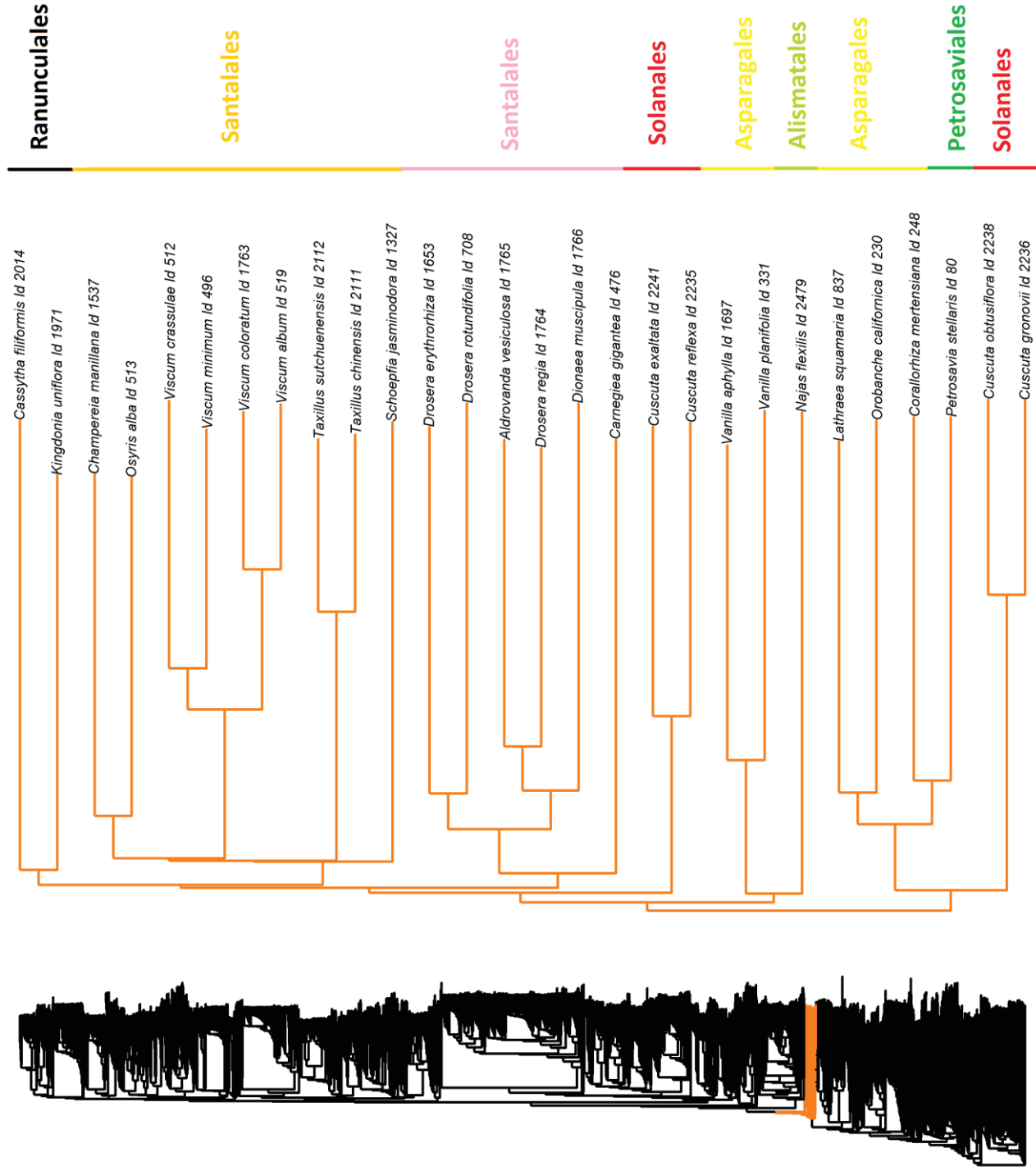


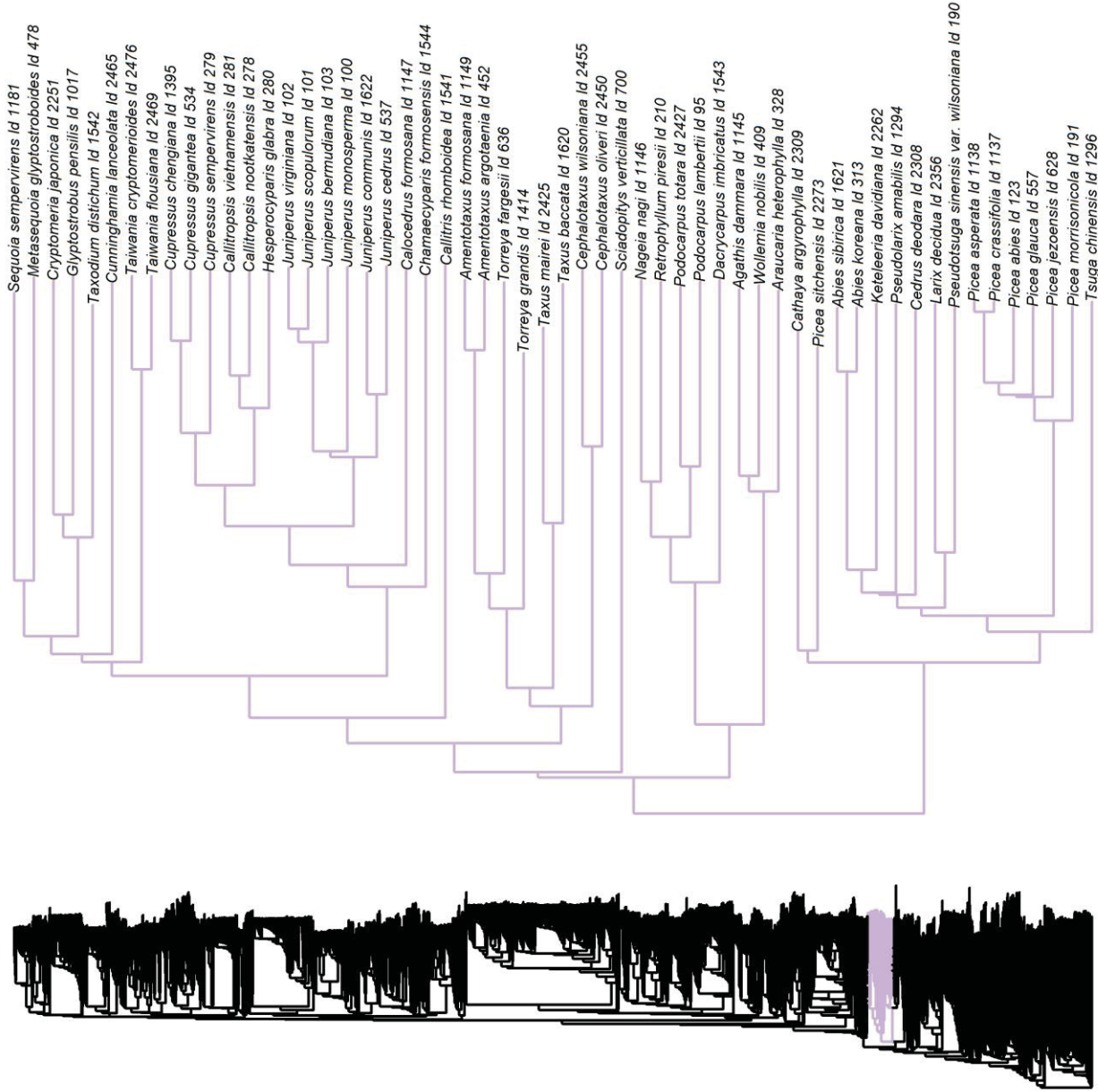


## Asterales



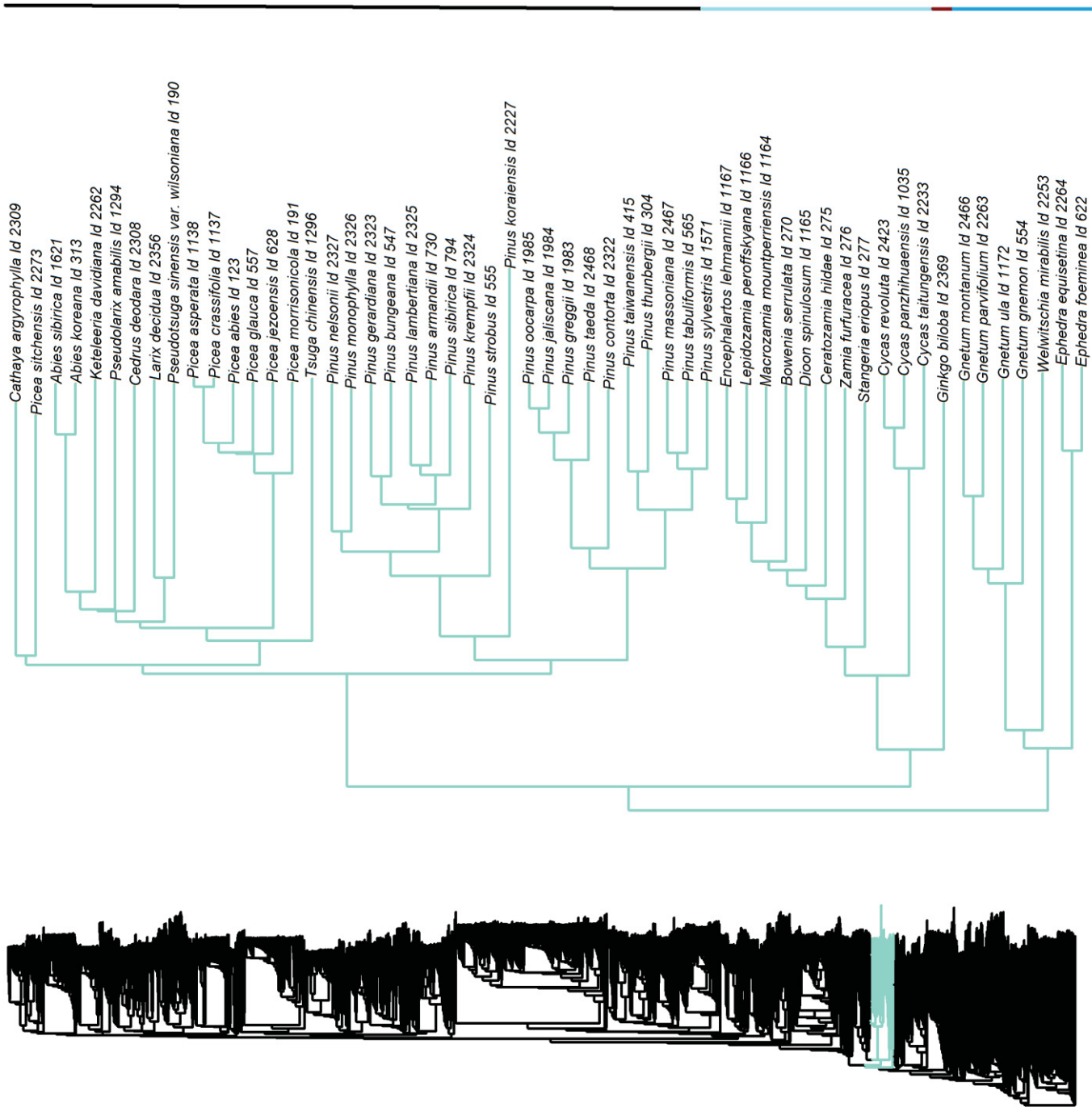
## Asterales

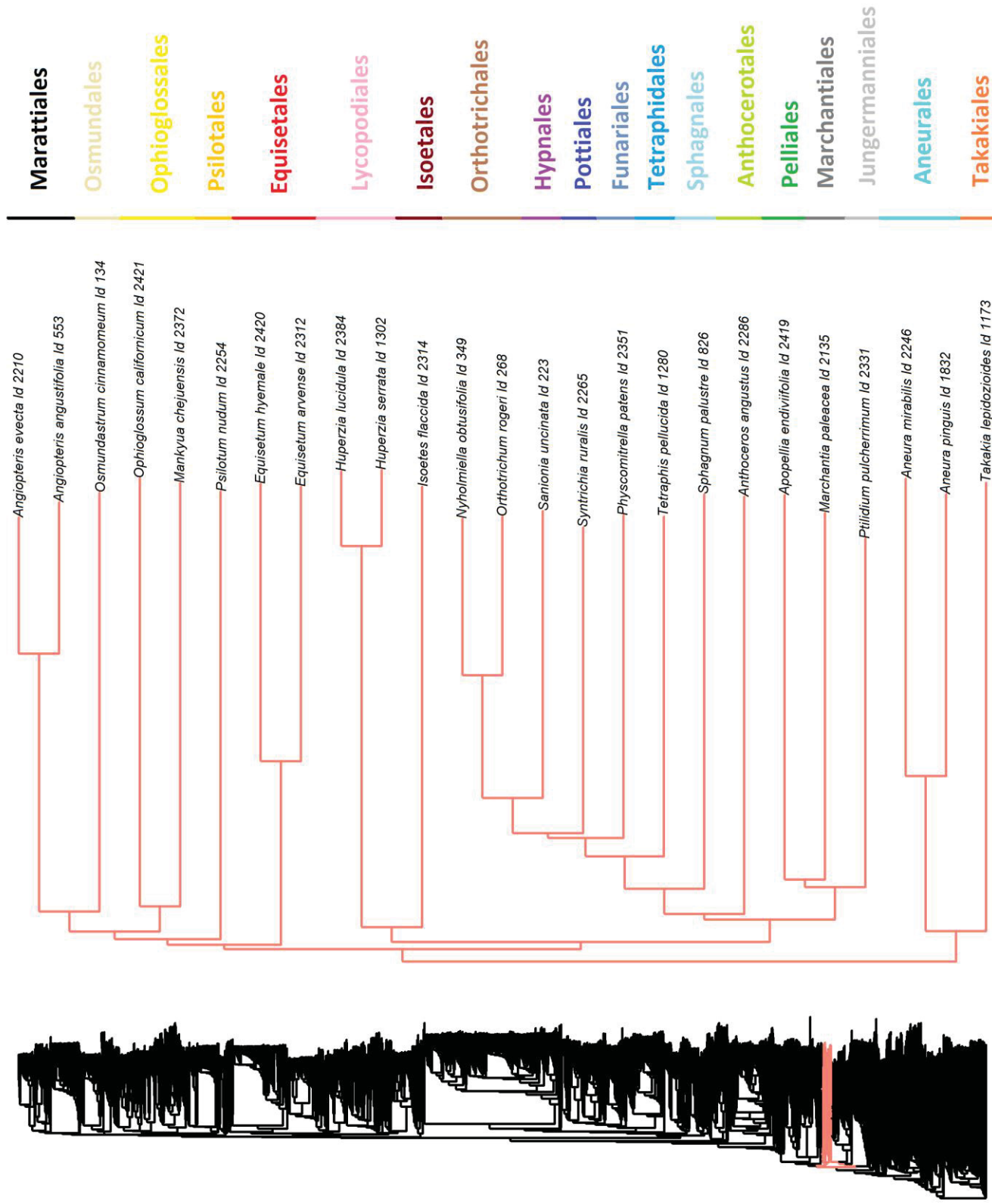


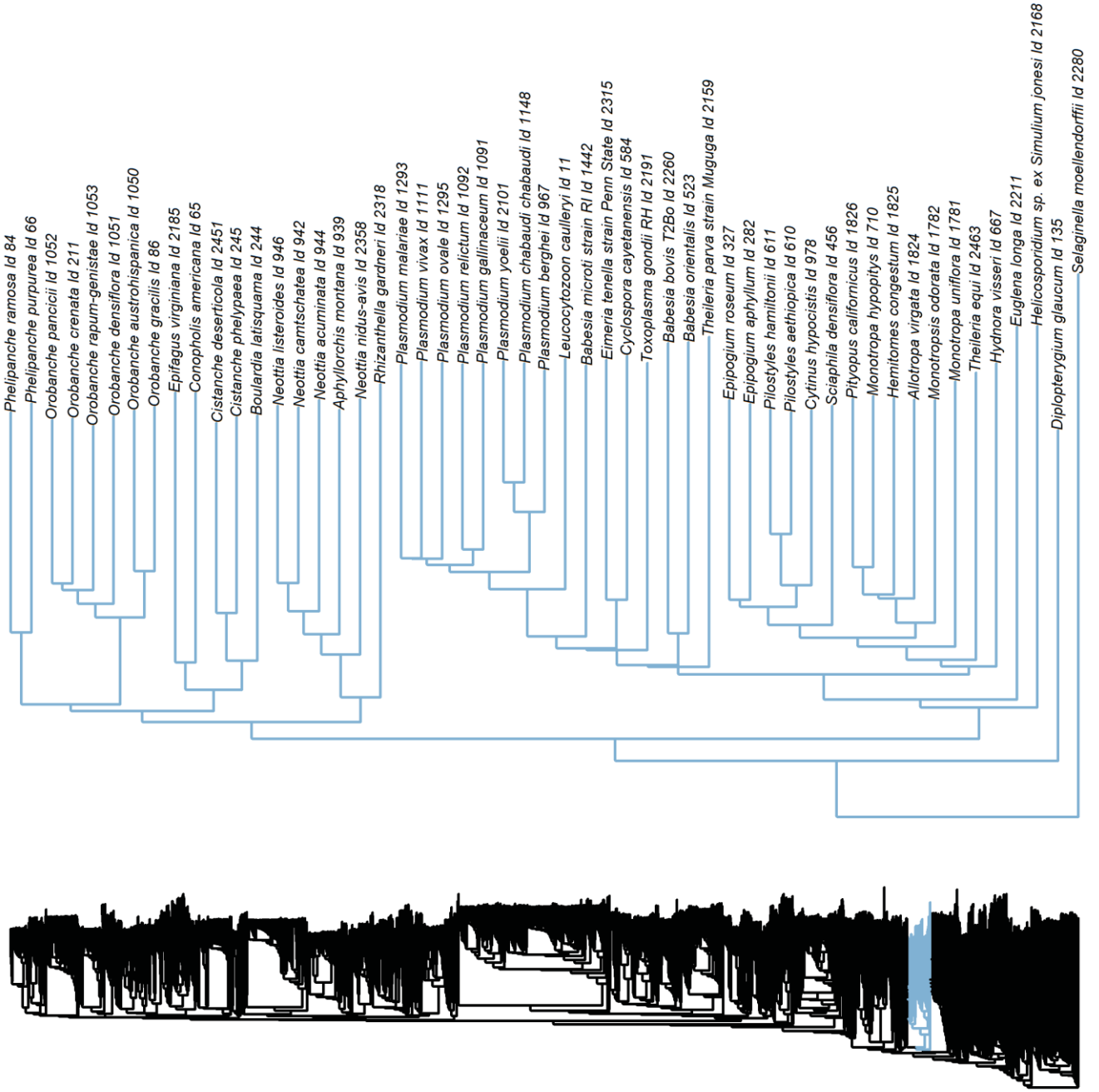


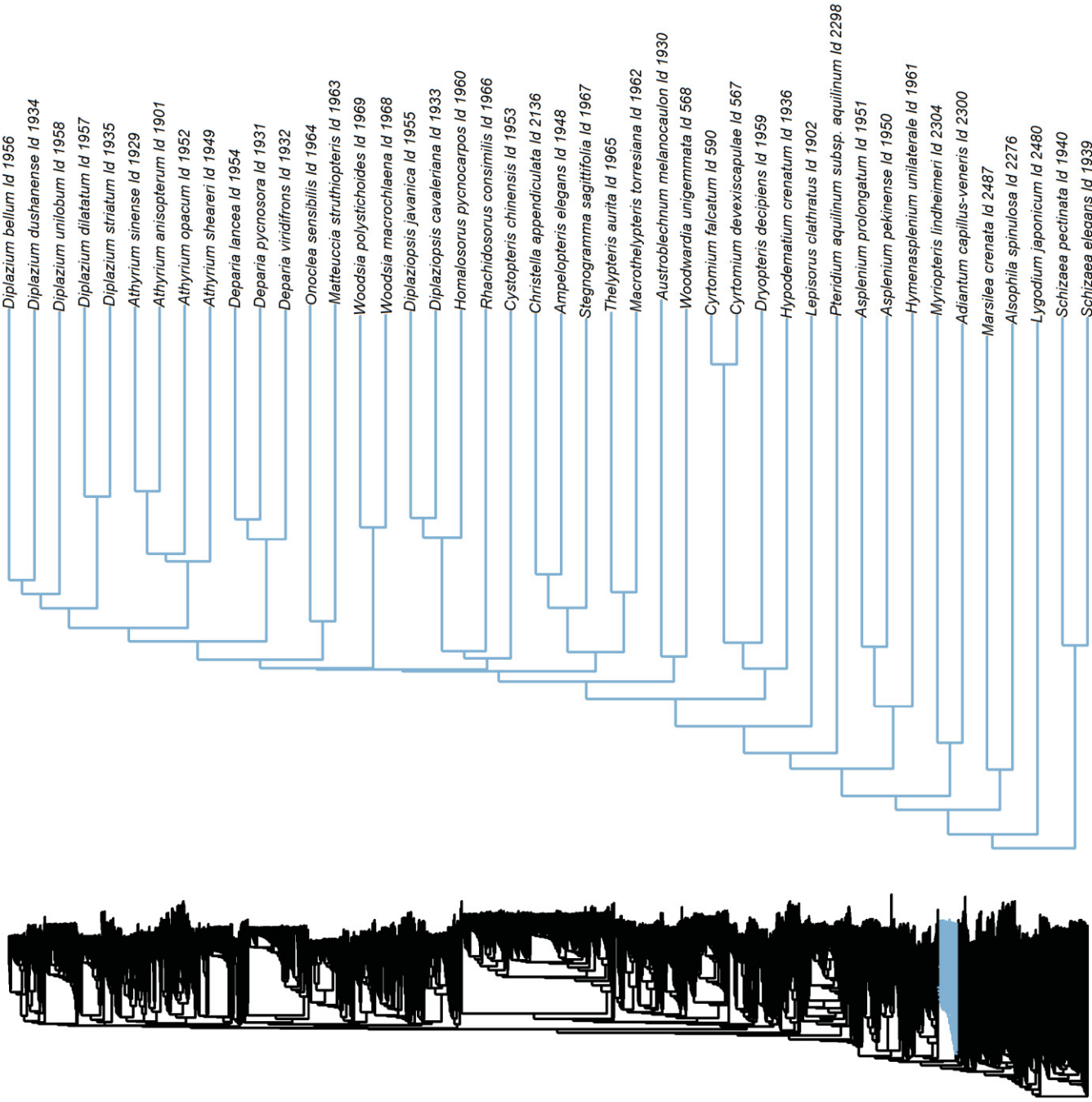
## Pinales





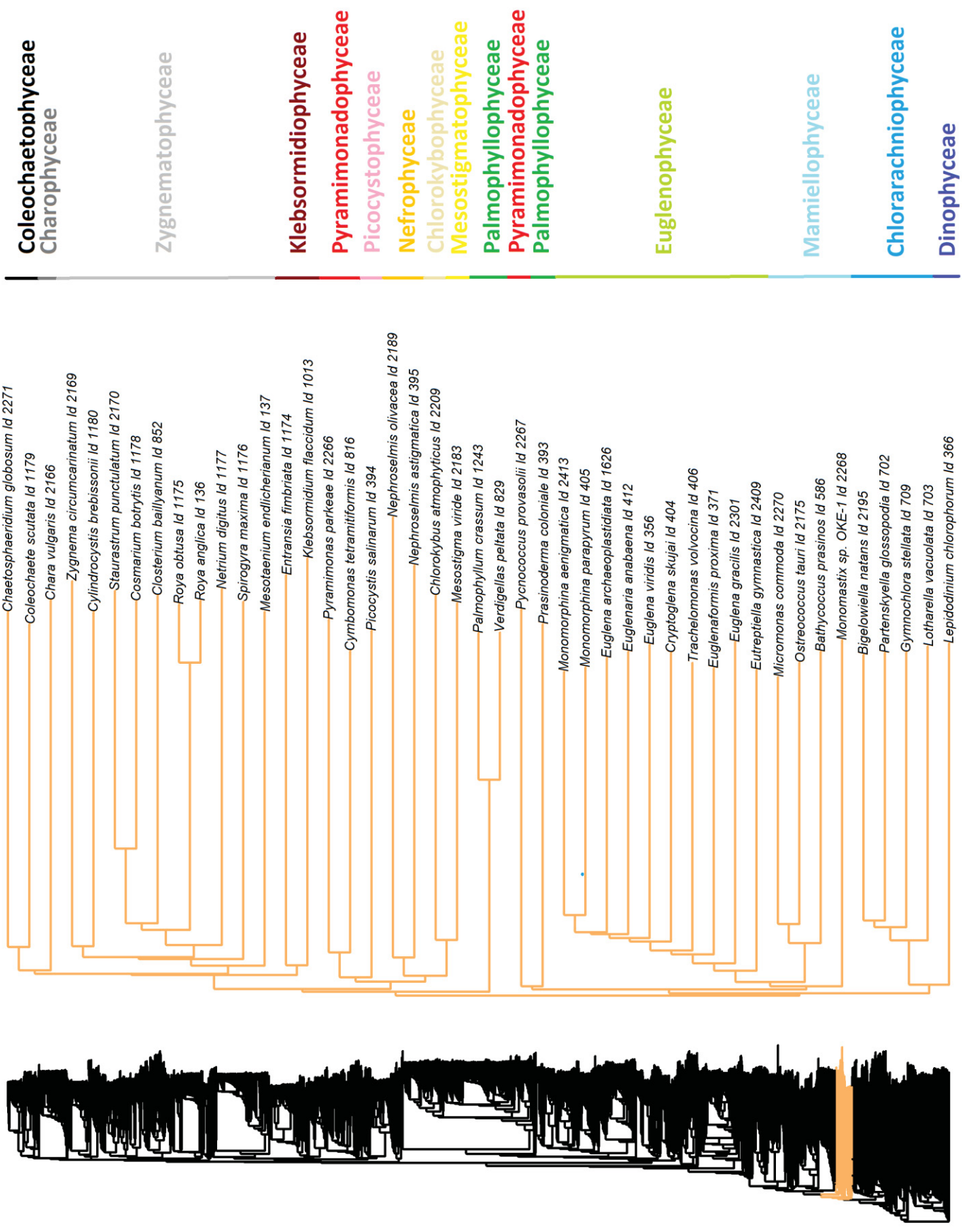




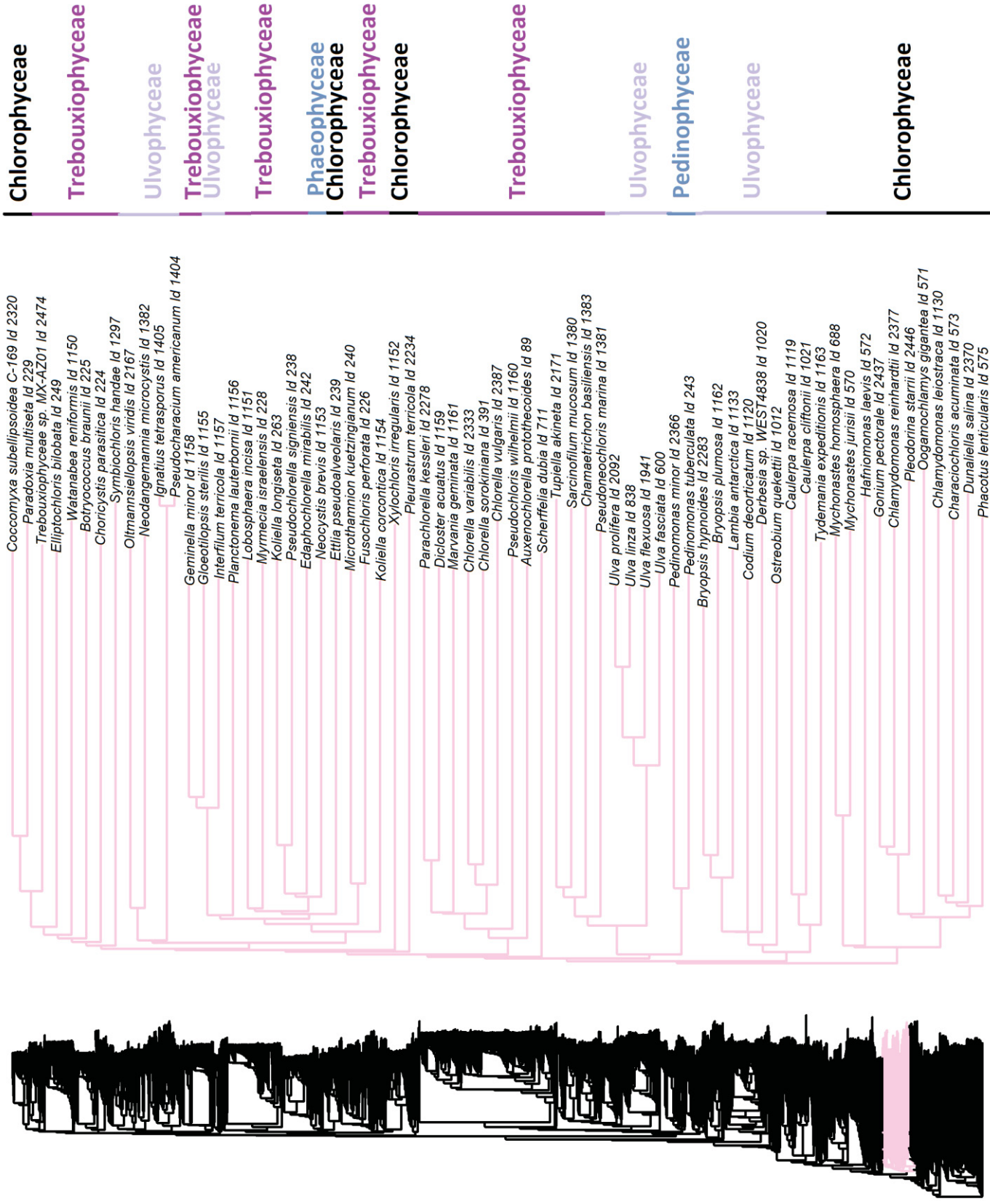


Polypodiales

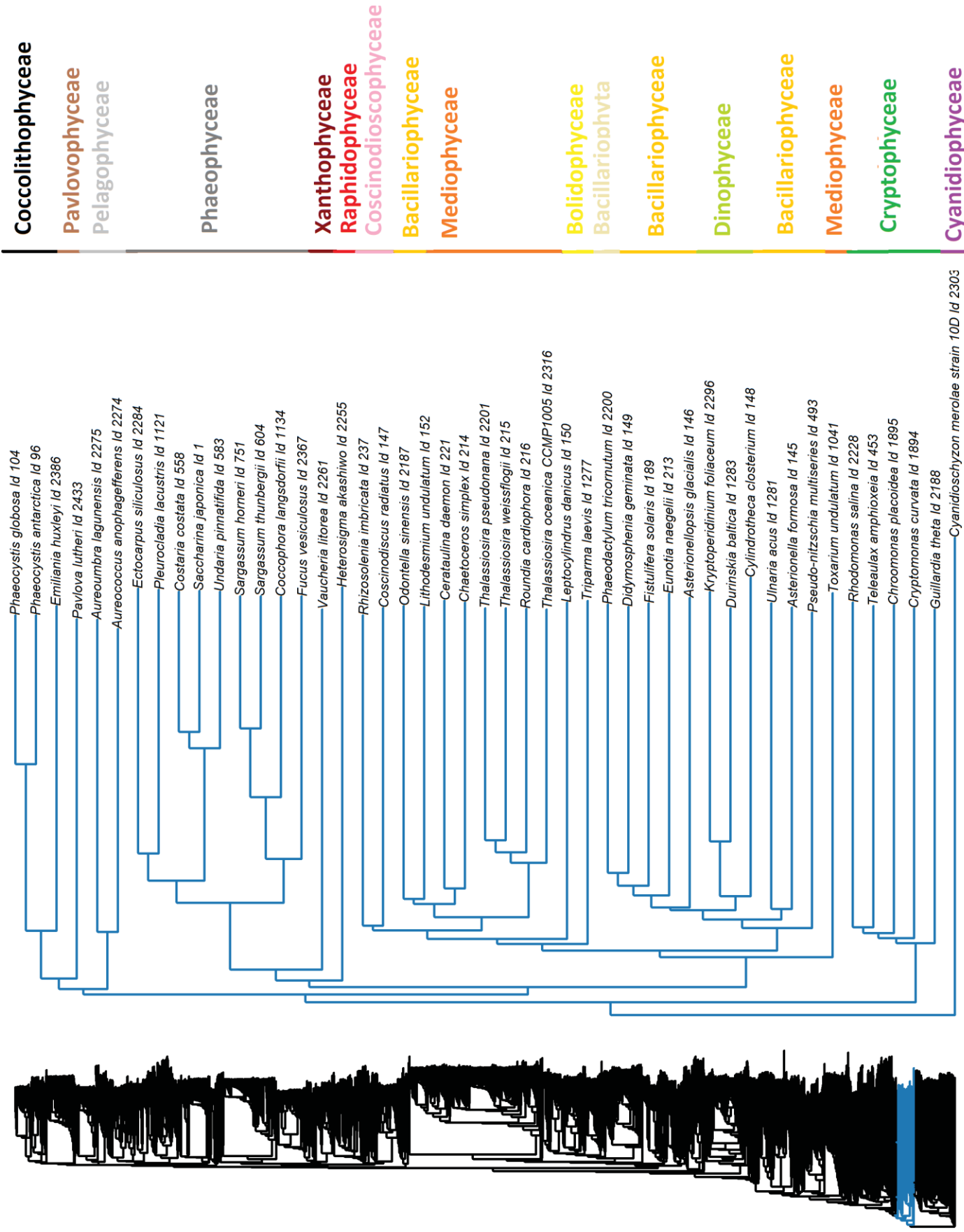
- Salviniales
- Cyatheales
- Schizaeales

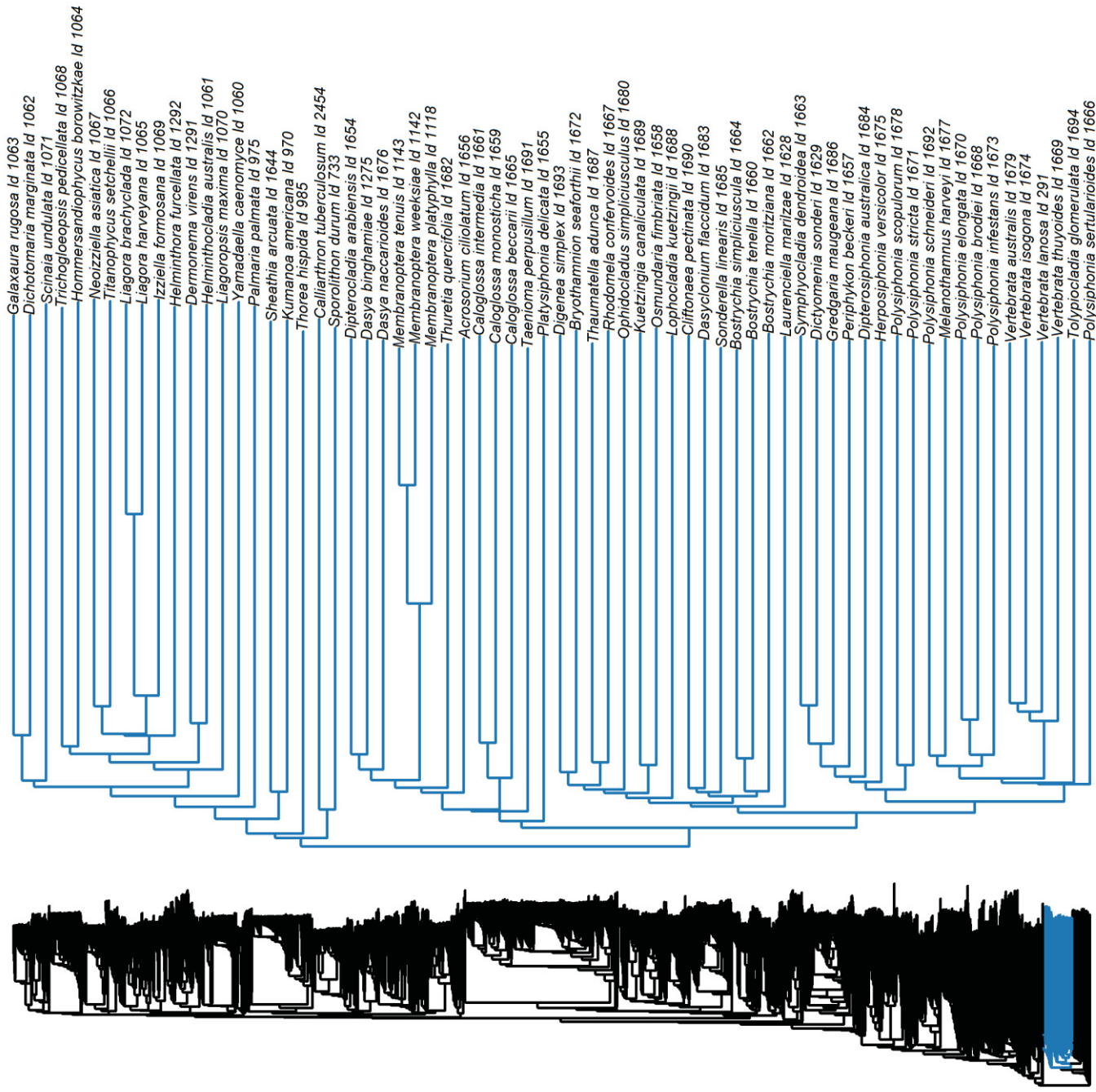




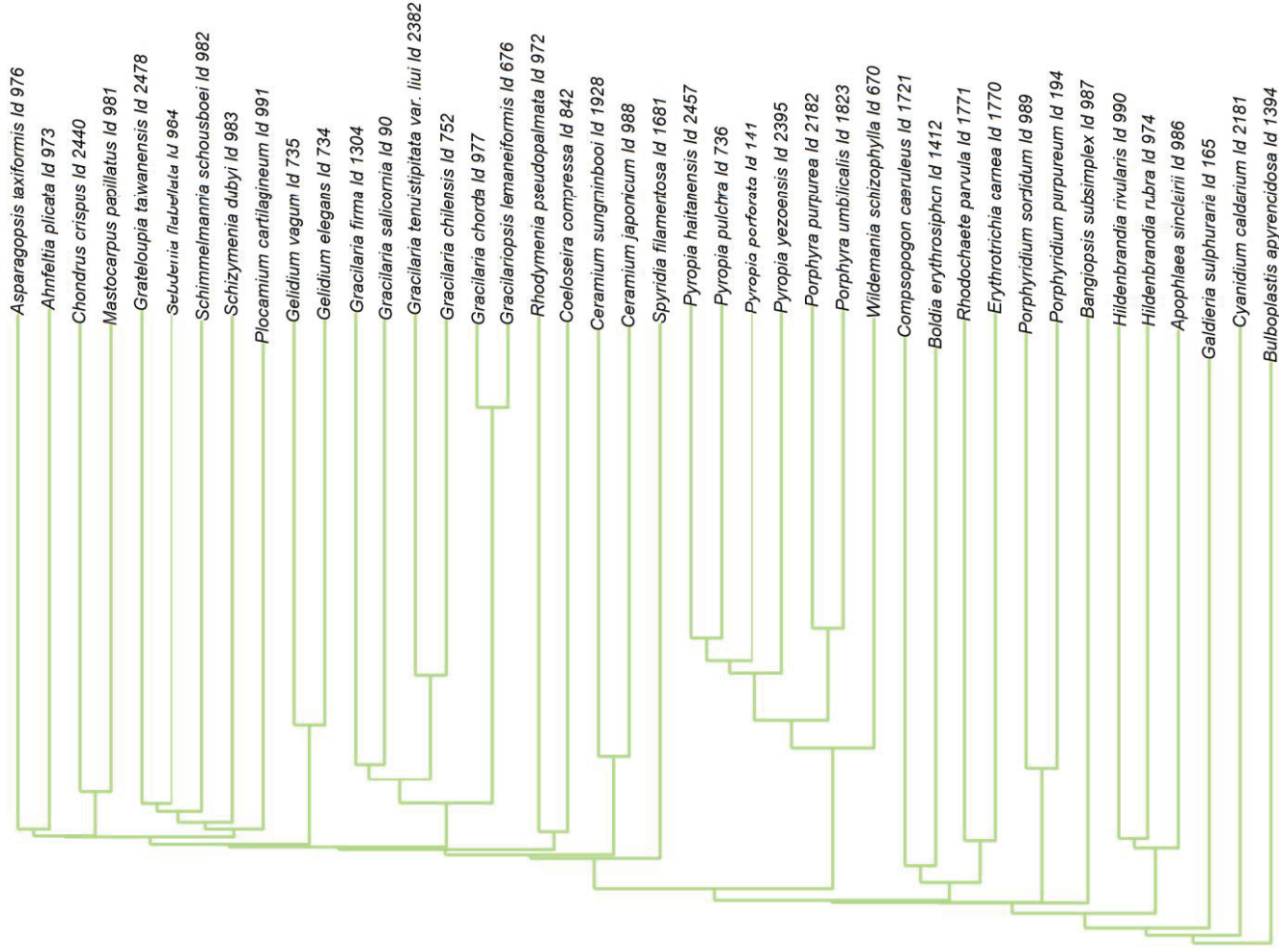








## Florideophyceae



Florideophyceae

Bangiophyceae

Compsopogonophyceae

Porphyridiophyceae

Stylonematophyceae

Florideophyceae

Cyanidiophyceae

Rhodellophyceae

## 6. ANEXO I: CERTIFICADO DE REGISTRO DE SOFTWARE SVECT




**REPÚBLICA FEDERATIVA DO BRASIL**  
Ministério Da Indústria, Comércio Exterior e Serviços  
Instituto Nacional da Propriedade Industrial

Diretoria de Patentes, Programas de Computador e Topografias de Circuitos Integrados

**Certificado de Registro de Programas de Computador**

**Processo nº: BR 51 2018 000244-7**

O Instituto Nacional da Propriedade Industrial expede o presente certificado de Registro de Programas de Computador, válido por 50 anos a partir de 1º de janeiro subsequente à data de Criação: 23 de novembro de 2016, em conformidade com o parágrafo 2º, artigo 2º da Lei Nº 9.609, de 19 de Fevereiro de 1998.

**Título:** **SVect - Slide Vector**

**Data de Criação:** 23 de novembro de 2016

**Titular(es):** UNIVERSIDADE FEDERAL DO PARANA

**Autor(es):** AMANDA WILCZEK  
/ ARYEL MARLUS REPULA DE OLIVEIRA  
/ BRUNO THIAGO DE LIMA NICHIO  
/ CAMILLA REGINATO DE PIERRI  
/ JERONIZA NUNES MARCHAUKOSKI  
/ JOSUE OLIVEIRA CAMARGO  
/ LETICIA GRAZIELA COSTA SANTOS  
/ MARIANE GONCALVES KULIK  
/ RICARDO VOYCEIK  
/ ROBERTO TADEU RAITTZ

**Linguagem:** MATLAB

**Campo de Aplicação:** 01, BL-07, IF-10, MT-01

**Tipo Programa:** IA-01, TC-01, UT-01

**Algoritmo Hash:** SHA-512

**Resumo Digital:**  
D54AAE745A1B064B83EE09D2431569DE96B996CE3AB32EC50B5CC221E513E99B6AF5127D8F7A43  
710DFD14B51648DBE6C29321EDFD4F167DB6040C8A38AD77E6

**Expedido em:** 06 de março de 2018

**Aprovado por** Julio Cesar Castelo Branco Reis Moreira